ARPANET COMPLETION REPORT

ORAFT

September 9, 1977

Isble of Contents

I. Executive Summary

I = 1

II. Formal Completion Report

 $II \rightarrow 1$

Profeca

- 1. Program Objective and Technical Need
- 1.1 Defense Program Addressed
- 1,2 State of the Art at Program Inception
- 1,3 Specific Technological Problem Addressed
- 1.4 Expected Payoff/Time Frame/Couts
- 2. Program Description and Evolution
- 2.1 Program Structure
- 2.2 Major Technical Problems and Approaches
- 2.3 Major Changes in Objectives and Approaches
- 3. Scientific and Technical Results and Accomplishments
- 3.1 Results of the Effort in Relation to the Program Objectives
- 3.2 Technical Aspects of the Effort which were Successful and Aspects of the Effort which did not meterialize as originally planned
- 4. Applications and Considerations for the Future
- 4.1 Conclusions of Technical Feasibility
- 4.2 Recommendations on Additional R&D Requirements and Opportunities
- 5. Program Impact and Assessment of Technology Developed
- 5.1 Service Use of Technology
- 5,2 Impact on Non-DoD Programs
- 5.3 Applications of the Program Results
- 5,4 Advance in the State=of=the=Art
- 6. Bibliography of Reports

7. Distribution List

III. A Review of the ARPANET Project

 $III \rightarrow I$

- 1. History
- 1,1 Background
- 1,1,1 Capsule History of ARPA and Its Information Processing Techniques Office
- 1.1.2 The State of the Art of Computer Communications -- Circa 1967
- 1,1,3 The RAND Study of Distributed Communications Networks
- 1.1.4 The Lincoln/SDC Experiment
- 1,1,5 The NPL Data Network
- 1.1.6 The Events Leading Directly to the ARPANET
- 1,2 The Events of 1967 and 1968
- 1.3 Key Aspects of the RFQ
- 1.4 Chronological Development, 1969 to 1975
- 1.4.1 The Groups and the Key People
- 1,4,1,1 Management and Administration of the Networks ARPA/IPT, DSS=W, RML, and DECCO
- 1.4.1.2 The Network Analysis Corporation
- 1.4.1.3 The Telephone Companies
- 1.4.1.4 Bolt Beranek and Newman Inc.
- 1,4,1,5 The Network Information Center
- 1.4.1.6 The Network Measurement Center
- 1.4.1.7 The Network Working Group and Host Involvement
- 1.4.2 Initial Subnet Design
- -1,4,3 Subnet Development
 - 1.4.4 Host Protocol Development
 - 1.4.4.1 Hest/Hest Protoce!
 - 1.4.4.2 TELNET
 - 1.4.4.3 Others
 - 1.4.5 Network Growth == A Summary

- 1.4.5.1 Traffic 1.4.5.2 Topology 1.4.5.3 Hosts
- 1,5 The Impact of the ARPANET 1,6 Maturity and Handover to DEA

Appendices

- 1.A BBN ARPANET BIDTIOGRAPHY
- 1.A.1 Technical Reports 1.A.2 Published Papers
- 1.8 Hest Protogo! Chronology
- 1.C ARPANET Logical Maps
- 2. Design and Implementation
- 2.1 Overview and Technical Basis
- 2.1.1 Terminology
- 2.1.2 Design Goals; Performance and functions
- 2.1.2.1 Basic Issues
- 2,1,2,2 Network Performance Goals
- 2,1,2,3 Intended Functions
- 2,1,3 Fundamental Design Choices
- 2.1.3.1 Transmission Facilities
- 2.1.3.2 Switching Technology
- 2,1,3,3 Topological Structure
- 2.1.3.4 Communications Equipment
- 2,1,3,5 Communications Protocols
- 2.1.3.6 Interface Characteristics
- 2.2 Topological Design
- 2.2.1 Introduction
- 2.2.2 Centralized Network Design
- 2.2.3 Backbone Network Design
- 2,2.4 General Network Design

- 2.3 Communications Equipment
- 2.3.1 The Network Circuits
- 2.3.2 The INP Hardware and Software
- 2.3.2.1 Network Considerations
- 2,3,2,2 The 516/316 IMP Hardware
- 2.3.2.3 The 516/316 IMP Software
- 2,3.2.3.1 General Descriptions
- 2.3.2.3.2 Data Structures
- 2.3.2.3.3 Packet Flow Through Major IMP Routines
- 2.3.2.3.4 Control Organization
- 2,3.3 The TIP Hardware and Software
- 2.3.3.1 The TIP Hardware
- 2.3.3.2 The TIP Software
- 2.3.3.3 The TIP Command Format
- 2.3.4 The NCC Hardware and Software
- 2.3.4.1 Introduction
- 2.3.4.2 The NCC Hardware
- 2.3.4.3 NCC Outputs
- 2.3.4.4 The NCC Softwere
- 2.3.5 The Pluribus IMP Hardware and Software
- 2.3.5.1 Hardware Structure of the Pluribus IMP
- 2.3.5.2 Software Structure of the Pluribus IMP
- 2.3.5.3 PLURIBUS IMP Reliability
- 2.3.6 Other Developments (SIMP, PLI, RJE)
- 2,4 Subnet Protocols (Intro, Overview, and Plan)
- 2.4.1 IMP+IMP
- 2.4.1.1 Basic Concepts
- 2.4.1.2 Original Implementation
- 2.4.1.3 Subsequent Modifications
- 2.4.1.4 Conclusions

- 2.4.2 The Routing Algorithm
- 2.4.2.1 Basic Concepts
- 2.4.2.1.1 NonmAdapt(ve Algorithms
- 2.4.2.1.2 Centralized Adaptive Algorithms
- 2,4,2,1,3 Isolated Adaptive Algorithms
- 2.4.2.1.4 Distributed Adeptive Algorithms
- 2.4.2.1.5 Routing Processing Goals
- 2,4,2,1,6 Routing Performance Measures
- 2.4.2.1.7 Routing Cost Measures
- 2.4.2.2 Original Implementation
- 2.4.2.3 Subsequent Modifications
- 2,4,3 Source IMP+Destination IMP
- 2.4.3.1 Baste Concepts
- 2.4.3.2 Original Implementation
- 2.4.3.2.1 Message Format
- 2.4.3.2.2 The Original Source=To=Destination Transmission Procedure
- 2.4.3.2.3 The Next Implementation
- 2.4.3.3 Subsequent Modifications
- 2.4.3.3.1 Dynamic Blocks
- 2.4.3.3.2 Restructured Message Numbers
- 2.4.4 HosteIMP
- 2,4,4,1 Basic Concepts
- 2.4.4.2 Original Implementation
- 2,4,4,3 Subsequent Modifications
- 2.4.4.4 Conclusions
- 2.5 Host Level Protocols (Overview, Layering)
- 2.5.1 Host-Host Protocol
- 2.5.1.1 Basic Concepts
- 2.5.1.2 Original Specification

```
2.5.1.2.1 Connection Establishment
  2.5.1.2.2 Connection Termination
  2.5.1.2.3 Flow Control
            Out-of-Band Signalling (Interrupt)
  2,5,1,2,4
  2.5.1.2.5 Other Functions
  2.5.1.2.6 NCP/Process Interface
  2,5,1,3 Subsequent Modifications
  2.5.1.4 Conclusions
  2,5,2 Initial Connection Protocol
  2.5.3 TELNET Protocol
  2,5,3,1 Besic Concepts
  2,5,3,2 Original Specification
            Network Virtual Terminal Organization
  2.5.3.2.1
  2.5.3.2.2 NVT Printer
  2.5.3.2.3 NVT Data Keyboard
  2.5.3.2.4 End#af#Line Convention
            Break Stanal
  2,5,3,2,5
  2,5,3,2,6
           Interrupt Signal
           Local Functions
  2.5.3.2.7
  2.5.3.2.8 Control Codes
  2,5,3,3 Subsequent Modifications
            Option Negotiation
  2,5,3,3,1
           The Expanded NVT
  2.5.3.3.2
           Revised Control Code Structure
  2.5.3.3.3
            Defined Options
  2,5,3,3,4
2.5.3.4 Conclusions
 2.5.4 File Transfer Protocol (FTP)
 2.5,4.1 Basic Concepts
 2.5.4.2 Original Specification
            Data Representation Types
  2,5,4,2,1
  2,5,4,2,2 File Structures
           Transmission Modes
 2.5.4.2.3
           Error Recovery
 2.5.4.2.4
 2.5.4.2.5 FTP Commands
```

- 2.5.4.3 Conclusions
- 2.5.5 Other High-Level Protocols
- 2.5.5.1 Remote Job Entry Protocol
- 2.5.5.2 Network Remote Job Service
- 2.5.5.3 Graphics
- 2.6 Parformance
- 2.6.1 Introduction
- 2.6.2 Theoretical Models, Measurements, Simulations
- 2.6.3 Delay
- 2.6.4 Throughput
- 2.6.4.1 IMP Bandwidth
- 2.6.4.2 TIP Bandwidth
- 2.6.4.3 Host Bandwidth
- 2.6.4,4 Pluribus IMP Bandwidth
- 2.6.4.5 Network Throughput Performance
- 2,6,5 Reliability
- 2.6.5.1 IMP Techniques
- 2.6.5.2 Pluribus IMP Techniques
- 2,6.5.3 Network Reliability
- 2.6.5.4 Lockups
- 3. Experience in Use and Operations
- 3.1 Server Systems
- 3.1.1 UCLA Campus Computing Network
- 3.1.2 TENEX Systems
- 3.1.3 Multics
- 3.1.4 ILLIAC IV
- 3.1.5 The Datecomputer
- 3.2 User Systems
- 3.2.1 The Terminal IMP
- 3.2.2 ANTS
- 3.2.3 ELF

- 3.3 Specialized Communities of Interest
- 3.3.1 Climate Dynamics
- 3,3,2 The Seismic Community
- 3.3.3 Secure ARPANET Use
- 3.4 Network Mail
- 3.5 Distributed Systems
- 3.5.1 The Resource Sharing Executive
- 3.5.2 The National Software Works
- 3.5.3 The TIP Service Facility
- 3.6 Networksbased Experiments and Research
- 3.6.1 AFCS Experiment
- 3.6.2 Speech Transmission
- 3.6.3 Packet Broadcast by Satellite
- 3.6.4 Internetting and Gateways
- 3.6.5 Endeto=End Security
- 3.7 ARPANET Operation and Maintenance
- 3.7.1 Control Functions in the IMP
- 3.7.2 Early Operation
- 3.7.3 The Network Control Center
- 3.7.4 IMP Software Maintenance
- 3.7.5 Hardware Maintenance
- 4. Observations
- 4.1 Social Issues
- 4.2 Technical Lessons
- 4,2,1 Terminal Handling
- 4.2.2 Reliability and Fault Isolation
- 4.2.3 Maintenance Management
- 4.3 Impact on Other Research Areas
- 4.3.1 Specialized Resources
- 4.3.2 Interprocess Communication
- 4.3.3 Personeto=Person Interactions

4.3.4 Conclusions

CHAPTER IL EXECUTIVE SUMMARY

<this chapter has not been written yet>

,				
			**	

CHAPTER II: FORMAL COMPLETION REPORT

Profese

In June 1968 an ARPA progrem plan entitled *Resource Sharing Computer Networks* was formally propored and approved. research carried out pursuant to this plan has since become known and internationally famous as the ARPANET. This ARPA program has greated he less than a revolution in computer technology and has been one of the most suggessful projects ever undertaken by ARPA. The program has initiated a ferwreaching effect on the Defense Department[†]s use of demouters as well as on the use of computers by the entire public and private sectors, both in the United States and around the warld. Just as the telephone. telegraph, and the printing press had fareseaching effects on human intercommunication, the widespreed utilization of computer networks which has been satalyzed by the ARPANET projest represents a similarly farwreaching shangs in the computers by mankind. The full impact of the technical set in motion by this project may not be understood for many YOAPS,

1. PROGRAM OBJECTIVE AND TECHNICAL NEED

1.1 Defense Program Addressed

The June 1966 ARPA program plan states the objective of the program in the following way:

The objective of this program is twefeld: (1) To develop techniques and obtain experience on interconnecting computers in such a way that a very broad diese of intercations are possible, and (2) To improve and increase computer research productivity through resource sharing. By establishing a network tying IPT's research centers together, both goels are achieved. In fact, the most efficient way to develop the techniques needed for an effective network is by involving the research talent at these centers in prototype estivity.

Just as timeshared computer systems have permitted groups of hundreds of individual users to share hardware and software resources with one another, naturely connecting dozens of such systems will permit

resource shering between theusands of users. Each system, by virtue of being timeshared, gan offer any of its services to enother semputer system on demand, the most important eriterien for the type of network interconnection desired is that any user or program on any of the networked semputers can utilize any program or subsystem evaluable on any other demputer without haying to modify the remote program.

This objective was an entirely new and different approach to an extremely serious problem which existed throughout both Defense Department and the acciety at large. The many shousands of computer genters in the Defense Department and the many other thousands of computer centers in the private and public sectors operate almost completely autonomously. Each computer center is forced to recreate all the meftwere and date files that it wishes to utiliza. In MARY seses this involves the complete reprogramming of software or reformatting of data files. duplication and redundant effort is extremely sortly and time consuming. The June 1948 ARPA progrem plem estimated that "such duplicative efforts more than double the national costs of creating and maintaining the activates. There had been other dompletely different attempts to address this problem, such as ettempts at lenguage standards for somputers, attempts at standardizing the types of herdwere, and attempts at automatic translation between computer lenguages. Although each such approach had some value and utility, the problems of trying to there computer software resources or files was truly enormous.

In addition to the general problem shared with the rest of the adjected dominity, the Defence Department glad faces sertein special problems having to do with trainings quoting the program plans

"Military personnel trained to use one manufacturer's equipment must often be trained again to use a different station to another. Machines procured from different manufacturers require as many different user training programs as there are machines thus inhibiting pesitive transfer of training that could accumulate through the rotation of military personnel. Those data files and programs which have common utility to many military organizations and installations must be stend, prested and maintained separately at each

different machine, Military systems interconnected in a distributed interestive network obviete such constraints."

Another objective noted in the program plan was to permit the linking of specialized computers to the many general purpose computer genters. The program plan noted this objective as follows:

with recent improvements in the hardware eres, is will become most cost effective to design and construct computers particularly efficient at specialized tasks (e.o., compiling, list processing and information retrievel). Making such machines available to all the somputer research establishments would significantly increase the complity of these other conters.*

This program was addressed at no less than changing the use of computers by the entire Defense Department. It was clearly intended that the use of such a computer network would permit resource sharing within and agrees the military services and throughout the Defense research community.

1.2 State of the Art at Program Incaption

By the date of the program plan in the late 1960s most of the apacific technologies required for a computer network had individually been achieved in some form, for example, there had been many connections of phone lines to computers (e.g., the SAGE system, Air Line reservations systems, and Time sharing systems). However, there had been only a very small number of stempts to connect computers together for the purpose of experimenting with the sharing of resources.

- In the early 1968s an attempt was made to link computers tegether at the Western Data Processing Center at UCLA for the purpose of enabling similar computers to perform load sharing. A similar experiment was also performed at Ball Laboratories and aphieved reasonable suggests for beveral years.
- A number of networks were constructed for the primary purpose of message handling, including a Westinghouse inventory central system and several air line reservations networks. The SITA network for air line reservations was supplicingly advanged in concept in the

mid 1968s, but details about \$ITA were generally not known in the U.S. computing semmunity. In any case, the techniques used for such message systems were special purpose in nature and were not readily transferable into the general area of inter-computer communications.

A direct progenitor of the ARPANET was an effort made in the mid 1968s to aghieve a soupling between assessing computing expertise and the operation of the SDC G32 demputer. This effort led to a phene line connection between the G32 at SDC and the TX2 at Lincoln Leboratory, and demonstrated the relative sees at modifying time sharing systems to permit network interactions.

Aside from the technical problems of interconnecting computers with communications circuits, the notion of computer networks had been considered in a number of please from a theoretical point of view. Of particular note was work done by Paul Barron and others at the Rand Carporetion in a study "Dn Distributed Communications" in the early 1966*s. Also of note was work done by Donald Davies and others at the National Physical Laboratory in England in the mid 1968*s.

In sum, at the time of the ARPA progrem plan in 1968 many of the requisite ideas had been gonsidered and many of the requisite technical bits and pieces had been attempted in some form, but no significant attempt had ever been made to put them together into a resource sharing computer network.

1.3 Specific Technological Problem Addressed

The technological problems of building the ARPANET gan be considered at many different levels of detail. At the top level, there were really two problems:

- In To construct a "subnetwork" considering of telephone circuits and switching nodes whose reliability, delay characteristics, capacity, and sest would facilitate resource sharing between the computers on the network,
- 2. To understand, design, and implement the protocols and procedures within the operating eystems of each connected computer, in order to allow the use of the new subnetwork by those computers in thering resources.

Within these two major technological problems, there were, of course, a lerge number of subsproblems; including the

engineering of the phone circuit connections, the topology of the network, the selection of switching node equipment, the design of line disgiplines to work through phone line errors, the routing problem, and many ethers,

1.4 Expedted Payoff/Time Frame/Costs

The goals for the ARPANET project were very broad and enviseded a significant eventual impact on the use of computers in both the public and private sectors. However in addition to these long range goals, ARPA visualized some quite specific initial payer in the form of improved productivity of the ARPA research program itself, and a resulting specimental benefit to the services from ARPA research. Guesting from the 1968 program plant

*Beventeen computer research denters throughout the country are supported in whole or in part by Information Processing Techniques. The installation of an effective network tying these location together should substantially reduce duplication and improve the transfer of scientific regults, as well as develop the network techniques needed by the military. The

research autput of these projects is important to all three Services and it is expected that this output can be substantially increased for the same dollar cost if a portion of the funds are utilized for the network."

In addition, initial payoff was antisipated in the form of technology transfer from the ARPANET project in three ways:

- a By dissemination of new scientific knowledge through conferences and the appropriate literature.
- o By transfer of management of the ARPANET to a gommon tarrier, and the resulting availability of ARPANET services to other groups (such as Office of Education Regions) Laboratories, NSF-supported universities, and various user groups supported by the NIH).
- o By adoption of the network technology by specific military groups (such as the National Military Command System Support Center and other military senters affiliated with its e.g., CINCPAC, CINCEUR, and MACY).

The expected events and the schedules for those events was stated in the June 1968 ARPA Program Plans

- a. July 1968 Award IMP contract
- b. March 1969 == Demonstrate (nitie) net operation with four nodes
- c. April 1969 Appreve design and extend sontract to include installation of 19 IMPs
- d. December 1969 == Complete network operational
- e. 1970 Add communication lines as necessary
- for 1971 = 1972 == Arrenge with a common cerrier the transfer of the communications system

The gosts enticipated by the Program Plan were as follows:

Yaar	Casta

Communication Line IMP Contractor Total

FY 68	6	262K	563K
FY 69	25K	1988K	1025K
FY 70	688K	\$00K	888K
FY 71	900K	199K	1000K

This sort estimate did not include the costs required for the various resource sharing experiments at the host sites in the ARPA community. The Progrem Plan avoided this issue, perhaps by necessity, as initials:

"The majority of this second siess of gosts will be bourned by each of the domputer research contrasts new extent. They will vary advess a wide range of extremes bounded by, for exemple, a single researcher's small experimental program and a group of researcher's centern with studies of oheline documentation."

2. PROGRAM DESCRIPTION AND EVOLUTION

2.1 Program Structure

With the signing of the program plan in June 1968, ARPA began work in earnest on three parellel paths of afforts (1) to obtain the network elequites (2) to select the system contractor for the switching nades and the everall designs and (3) to initiate efforts within the ARPA research community for resource sharing experiments and specialized network supports

The ARPA IPT office had not had great experience in acquiring communications circuits and was rather piecently supprised to discover that enother DeD organization, DECCO, was able to handle all the contractual details with the common serviers for circuit lesses. Most of the required 50 kilobit circuits used in the ARPANET were lessed through DECCO from ATET, but a small number of circuits were lessed from other carriers such as General Telephones. In addition, ARPA was fortunate in being able to arrange a rather high-level contact in ATET (long lines), which greatly facilitated the interactions between the network system contractor, ARPA, and ATET, The selection of network hade locations and the internode connections (and,

therefore, the location of circuit terminations) was a specialized topology problem and represented a difficult theoretical problem in its own right. To help solve this particular problem, ARPA contrasted with the Network Analysis Corporation (NAC). NAC had developed gertain networking analysis tools and via this ARPA support, such tools were refined; NAC's advice on topology was sought through the various stages of ARPANET growth.

In a progedure relatively unusual for the IPT affice in ARPA, a competitive procurement was planned for the selection of the contractor to design the switching nodes and set as general systems contractor. An RFP was proposed and (saued on (), and efter an evaluation of approximately 15 bidders, Selt Beranek and Newman Ing. (BBN) was selected as the systems contractor for the degion of the subnetwork and the switching nedes: BBN subcontracted with Honeywell for the switching node hardware tacks. Over the 140e of the ARPANET program from January 1969 until the transfer to DCA in July of 1975, BBN served as the eyetome contractor for the design, implementation, operation and maintenance of the subnetwork, and in 1977, BBN is still serving that role under the aegis of the Defense Communications Adency (DCA).

In a samewhat less structured way, the research groups receiving ARPA IPTO support were then encouraged to begin considering the design and implementation of pretegols and procedures and, in turn, computer progrem modifications, in the various host computers in order to use the subnetwork. specific responsibilities were arranged; UCLA was specifically asked to take on the task of a "Network Measurement Center" with the objective of atudying the performence of the network as it was built, grawn, and medified; \$RI was apecifically asked to take on the tack of a "Network Information Center" with the objective of golfectine information about the network, about host resources, and at the same time generating samputer based tools for storing and accessing that collected information. these two specific contracts, some rather ad hat mechanisms were pursued to reach agreement between the VARIOUS research the appropriate "host eant regtors about protocols" for intercommunicating over the subnetwork. The "Network Working Group* of interested individuals from the various host sites was rather informally angulaged by ARPA, After a time, this Network Working Group became the forum for, and eventually a semi-official approval authority for, the discussion of and issuance of host protocols to be implemented by the various research contractors. Since this pert of the program was much less structured, programs was rather slow for a time, but with some ARPA IPTO pressure on the research contractors, this mechanism eventually was surprisingly successful in establishing effective host protocols.

2.2 Hejor Teghnigal Problems and Approaches

In a program of this duration and complexity, it is not difficult to identify many dozens of important technical problems and approaches to those problems. We have list a few of the problems which were most technically challenging in the early few years of the ARPANET program. A few additional major technical problems will be listed in the next section on "Mejor Changes and Objectives".

TOPOLOGY - For any network of this type with even a degen nodes, an obvious, seriy recognised, and quite formidable problem is topological optimization.

Assuming that the node locations are known, the number of ways of arrenging M links among N nodes is very larger the links are usually evaluable in discreet sizes.

(bandwidth); and the complement cost structures, time delay functions, and reliability functions are typically non-linear. The design may be subject to many constraints, including maximum or average time delay, everage or seak throughout requirements, and reliability requirements. The usual deal of the optimization is to provide a network design that meets all constraints at the lowest cost. The approach to this problem was to design an elaborate computer program to assist in the optimization. It is not possible to use an exhaustive appreach, and instead the approach used was to demorate a "starting network" and them to perform toest transformations to the templogy in order to reach a locally optimum network. If this precedure is repeated with many starting networks and if the resulting locally optimum networks are evaluated, it is possible to find a feasible solution with costs that are aless to optimels

D LINE ERRORS - A critical necessity for a resource sharing gemputer network was to provide reliable sommunication and one component of such reliability was an ability to work through the expected phone line

errors on the 58 kilobit sircuits. The approach taken was to design apacial checksumming hardware at the transmitting and receiving and of each 50 kilobit singuit in the network. As part of the switching node transmission precedure, a powerful 24 bit cheaksum is appended to every packet of information to be transmitted. Then, at the receiving mode, the sheeksum is recomputed in hardware and compared with the transmitted sheeksum. If there is no error, an seknowledgment (with its own checksum) is sent back to the transmitting mode, If there ta an APPOF. acknowledgment is sent and the packet w111 **D** • retrensmitted.

design of an interface between the switching node and heat computers of many different types. It is important to allow a logical metal between the switching node sommuser word length and the verying word lengths of the heat computers, and also to allow the input/output routines of the switching node and the input/output routines of the heat computer to service the interface

gooperatively without placing an undue processing burden or tight timing constraints on either machine, The approach taken was to design a bit-serial interface which could legisally step after any bit, and to then require that each host computer build (obtain) a specialized small herdwere unit called a *special host interface* which would be installed in the host machine itself to serve the network connection. Thus, a locical and effectivel standard was specified by the switching network in a manner intended to minimize everall trouble to all hosts, then each host was required to meet that standard both in hardwere and settware.

SWITCHING NODE PERFORMANCE = A sentral meal of the network was to provide resource sharing between remote installations in such a way that a local user could employ a remote resource without degredation. In particular, this required that the subnetwork be sufficiently reliable, have sufficiently low delay, have sufficiently great capacity and have sufficiently low cost, that remote use would be "as attractive" as local use. These objectives translated into a major technical problem for the switching node itself to provide low deley and high separity at modest sost. The approach taken was to select a mini temputer for the switching nodes whose I/O system was very efficient and to write very carefully tellered programs in maghine language to optimize the especity and the lew delay of the data path in the switching node. Great extension was paid to minimizing the especialism time of the inter loops of these programs.

REMOTE CONTROL - A speaker and new problem was posed by the need to put dozens of small identical mini computers in the field, in an environment where the host connections to those computers was comewhat experimental and where the precreams in the mini computers themselves were under development and were changing from month to menth. It was important to be able to do debugging of software, debugging of new heat connections, program modifications, and the installation of new programs without the casts and difficulties associated with having menned sites. The approach taken was to develop

an entirely new technology of remete computer management. A Netwerk Centrel Genger was established for the subnetwork and software was developed which made it passible to examine or change the operating software in any node of the net from the central network control denter. This approach made it passible to issue complete new versions of the software to each node in the netwerk from a sentral place in an hour or two. In addition, each node reports the state of its health to the central place periodically and provides information on which to been debugging and maintenance activities.

ROUTING = In a non-fully connected network, an important problem is the decision process by which each node decides how to route information to reach any particular destination. This is a difficult theoretical problem and there are many different approaches, including fixed routing, rendem routing, centrally controlled routing, and various forms of distributed adaptive routing. The approach taken was to use a distributed adaptive traffic routing algorithm which estimates on the basis of information from adjacent nodes the globally proper

instantaneous path for each message in the face of varying input traffic loads and local line or node failures. Each IMP keeps tables containing an estimate of the output circuit corresponding to the minimum delay path to each other IMP, and a serresponding estimate of the delay. Pariodically, each IMP sends its current routing estimates to its neighbors; whenever an IMP receives such a message it updates its internal estimates. Thus, information about shanging conditions is requirely propagated throughout the network, and whenever a packet of traffic must be placed on the queue for an output gircuit, the imp uses its latest estimate of the best path.

HOST PROTOCOL - In many ways a computer network is to host computers as the telephone system is to human users: a transparent communication medium in which even after the caller has learned how to insert dimes and disi, it is etil! necessary that he speak the same language as the person called in order for useful dommunication to occur. The common language is referred

to as host protocol, and the preblem is to design a host protocol which is sufficiently powerful for the kinds of demandation that will occur and yet can be implemented in all of the various different host demputer systems. The initial approach taken involved an entity dailed a MNetwork Control Program which would typically reside in the executive of a host, such that processes within a host would dommunisate with the network through this Network Control Program. The primary function of the NCP is to astablish gennections, brack connections, switch demnections, and gentrol flow. A "layered" approach was taken such that more complex procedures (such as File Transfer Procedures) were built on top of aimpler procedures in the host Network Control Program.

2.3 Major Changes in Objectives and Approaches

The ARPANET development was an extremely intense activity in which contributions were made by many of the best computer scientists in the United States. Thus, almost all of the "major technical problems" already mentioned received sontinuing attention and the detailed approach to those problems changed

several times during the early years of the ARPANET effort. However, in addition several more major changes in objectives were introduced once the initial network became operationals

THE TIP - The initial model switching units, colled 0 IMPs, were intended to interconnect semputers and high bandwidth phone lines. At the outgot all terminal accous to the natwork was via terminal commentions to the heat themselves. After a time it became clear that there was a population of users for which terminal aggess to the network was very destrable, but who were not conveniently able to access the network via a host computer. Thus, a new model switching unit, a Terminal Interface Message processor, or TIP, was defined to serve the purpose of an IMP plus an additional function of direct terminel access. This shift resulted in the dealen of a TIP which really was a tiny heat embedded in a switching node itself and permitting the direct tennection of up to 63 asymphronous characterseriented terminels to the switching node. The TIP became the node! switching unit of shoics, often even where there was a issel host computer; this allowed connection of both hosts and terminals at that location directly to the network.

- PLURIBUS As the network grew the projected network traffic and the projected famout of hosts and terminals at a node began to expeed the capacity of the IMPs and TIPs in the Honeywell series. In order to plan for a future in which there might be sirguits of I megabit or more, it was decided to embark on the canatruction of a new switching nade which could handle higher bandwidths and higher fenerals. The approach taken was to dealign a multiprocessor switching nade which would parmit modular growth as the bendwidth and feneral requirements increased. This approach represented an interesting and important departure in computer technology as well as a necessary approach to achieving greater performance in necessary approach to achieving greater performance in necessary approach to achieving greater performance in
 - o SIMP The initial ARPANET directs were all 58 kilobit land lines. As the program developed, it became desirable to consider the use of satellite links and further desirable to attempt a more sophisticated use of

3. SCIENTIFIC AND TECHNICAL RESULTS AND ACCOMPLISHMENTS

3.1 Results of the Effort in Relation to the Program Objectives

In some cases a major program can be seen to reach its objectives in a single instant; a mushroom cloud at Alamogordo, or Armstrong atepping on the meen. In other cases like the ARPANET, an equally important objective may be reached with equal success, but the event must be observed in a more somplicated way and over a longer period of time. The first ARPANET objective "to develop techniques and obtain experience interconnecting computers in such a way that a very bread class of interactions are possible. Not just techniques, but an entire technology of packet switching has been developed, and an enormous body of experience has been developed on intersemmenting computers to allew a broad diess of interestions. The second objective was "to improve and ingresse domputer research productivity through ressures sherings. Computer productivity has been improved and has been increased through resource sharing over the ARPANET. Further, computer research productivity has been improved and intreesed in even some unexpected ways by an ARPANET-induced sealed change; within the

DRAFT Formal Completion Report

ARPA research community, it has become more <u>factionship</u> to work placely with other remote groups of researchers.

Another objective was to permit the linking of specialized computers to the many general purpose somputer senters. Several mejor specialized computers have been linked over the ARPANET to the main general purpose semanter senters in an extremely successful fashion.

At a technological level, the overell objectives were to construct a subnetwork whose reliability, delay sharecteristics, deposity, and cost would fasilitate resource sharing between the computers in the network; and to understand, design and implement the Host protocols to permit such resource sharing. Such a subnetwork, consisting of over 38 nodes and stretching from Hawaii to Norway, was suggessfully constructed; and the necessary Host protocols to ellow resource sharing between the connected computers were successfully understood, designed and implemented.

From a programmatic point of view, the initial IMP contract was awarded a months later than had been anticipated by the program plan, but the demonstration of initial net operation with a nodes and later the extension to 19 modes took piece on a time

The first, salled "network meil", was a case relatively simple idea which, in rudimentary form, had been a "curlosity" at some installations for years, suddenly and quite unexpectedly became very successful. The surrent form of the is called "network mail", and provides computer "mailboxes" for many individuals in their favorite Host computer. Whenever A desires to communisate a message to B, C, D, and E, he logs into a convenient Host, uses a program that sends mail, lists the addresses with an indication of their mailboxes, and types the message. No matter where B. C. D. and E might be, they occasionally enter the net at a local nede, too into their (mailbox) Host, and as part of the login process sutomatically told that a message exists for them. They may them obtain any recent messages by using another program to read their mail. The desire to provide mail feelilities of this type became an end in itself and partially influenced network growth and modification.

The second, use of the ARPANET for digitized speech experiments, was a seen where improvements in the vocader technology (and the continuing defense need for encrypted speech) led to a serious investigation of speech transmission over packet

networks. Singe speech traffic posed different requirement on network performance (e.g., low delay) changes were required in the network IMPs to accommodate this traffic. This activity is still growing in importance and is now being pursued by ARPA separately from the ARPANET program.

the satellite channel than simple point=to-point connections. In particular, the approach taken was to consider multi-escess her of a single up frequency and a single down frequency where all stations den hear all other transmissions including their ewn. This approach required the development of a new kind of switching node called a Satellite Interface Message Presessor (SIMP) which dould hendle the complex algorithms required for multi-eccess use of the satellite channel. In the final years of the ARPANET project prior to the transfer of the ARPANET to DCA, the initial work on such Satellite IMPs was accomplished. However, the full utilization of this squipment is still in progress under ARPA support and is being pursued separately from the ARPANET project itself.

PLI - The ARPANET was initially intended for use in an entirely unclassified fashion between ARPANET-supported research groups. However, a gleep intention of the segivity was to understand techniques which might be applied to defense department needs, and one such need implied a desire to use the ARPANET for classified

traffig. Since all of the circuits and most of the nodes existed on entirely unclassified terrain, an and-tamend eneryption scheme was needed. However, this seeed technical problems because all network nodes regular addresses in the clear, and because no evallable engryption units metahed the formats required by the ARPANET IMPs and TIPs. The approach taken was to design assemblace of two minicomputers and a KG=34 enerystion unit to be a "Private Line Interfage", or PLI, and to allow elegatified traffic between two secure hosta to take place on a HostuPileARPANET-PileHost connection. This project was conducted in cooperation with NSA, and had the additional interesting attribute that it required (probably for the #frat time) the "sertification" by NSA of minicomputer software to be used as part of eneryption achome.

In the letter years of the ARPANET project, two "applications" of the ARPANET become sufficiently important to influence the overall objectives of the program,

again that was extremely gloss to the (ngramental time anticipated in the program plan. However, as the success of the ARPANET became abvious after about 2 years, decisions were made to grow the net to a size that had never been anticipated in the program plans growth to more than 58 nades ever a desgraphic area from Hawaii to Norway was dertainly not priginally anticipated. In similar feshion, the seat in the first 2 fiscal years of the program were very close to enticipated costs, but as decisions were made to expand the scale of the network, the cost no longer followed the program plan.

In short, the results of the effort were eminently successful and far more than adequately met the objectives initially stated. The success of the program far exceeded even the most optimistic views at the time of incestion.

3,2 Technical Aspects of the Effort Whigh was Successful and Aspects of the Effort Whigh Did Not Meterialize as Originally Enviseded

Sings the ARPANET progress as a whele was a major suscess and sings the whole really is the sum of the parts, it is given that

the many key technical aspects of the program were in turn successful. A modest list of technical successes is as follows:

- o pewerful computer-based techniques for topological optimization were developed and the choice of ARPANET topology was made with the sid of these tools;
- o the sepriors successfully provided high reliability SOK/sec circuits;
- a subnetwork instuding nodel switching IMPs and TIPs was constructed whose performance, reliability, and cost did facilitate resource thanks;
- The ARPANET provides a convincing demonstration that adaptive routing algorithms can be made to perform reliably (e.g., in a globally correct manner in the face of local failures), efficiently (e.g., adapting to thanges in the network quickly and accurately), and flexibly (e.g., accommodating a veriety of circuit bandwidths and intermede distances) without excessive complexity and everhead.

- The ARPANET has demonstrated that it is possible to build a large operational natwork in such a way that the sifects of component failures are localized rother than "creshing" or otherwise making nonesperational large portions or all of the network. A nade or a Host can fail in the ARPANET and network use will be prevented for only the few users directly connected to that node or using that Host.
- o The ARPANET has confirmed the theoretical result that networks which store-endwionward packets gan aghieve delays which are low when tempared to the delays incurred in the computers (Hosts) which are uping the networks
- o The ARPANET has demonstrated that a network can be constructed so that nodes, lines, traffig, and so on can be added at deleted without major upheavals with each addition.
- o The ARPANET has demonstrated that it is passible for a network to control and operate itself without explicit control from a control center.

- The ARPANET has demanstrated new techniques for menitering, maintenence, and debugging. The nodes in the ARPANET typically operate at sites where there is no knowledgeable person evaluable locally. The nodes automatically report their status (via the network itself) to a network monitoring center, and maintenance and debugging (of both the software and hardware) are typically derpied out (again via the network) from the monitoring center.
- o Possibly the most difficult test undertaken in the development of the ARPANET was the attempt we which proved successful we to make a number of independent Host computer systems of varying manufacture, and varying operating systems within a single manufactured type, communicate with each other despite their diverse sharesteristics.
- o the Pluribus multiprocessor technology, important in its own right, was developed to provide an improved performance nodel switching units

o a set of Host protocols was hammered out between the Host organizations and resulted in a layered structure of Host protocols that did facilitate resource sharing.

Inevitability there were technique areas which proved less tractable than had been hoped and some kinds of network utilization developed more slowly than had been entidipated.

Problems of reusing, flew control, and congestion control in the subnetwork turned out to be more difficult than had been originally expected. These signrithms required intensive investigation and were medified several times in the early years of the network's growth. Luckily, the improvements in the signrithms managed to stay slightly shead of the growth in network size and traffic and, therefore, difficulties with the signrithms never represented an impediment to resource sharing on the network.

It proved more difficult to agree on the specification of adequate Heat protection then had been originally expected and this difficulty deleved wide use of the network for at least a year langer than had been enticipated. However, the eventual design of the Heat protects were eminently

successful and provided a strong base for resource sharings

of a widely-useful on-line Network The development Information Center (NIC) proved to be a difficult project. Although very interesting and important computer=based tools For information handling were developed, the actual timely collection and on-line publication of detailed deteriotions of resources at the various Host sites proved meanly intractable. Further, the somputer toois that developed while rather elegant were not so elearly cost effective for wide socie remote use. Eventually, the socia of the NIC were substantially reduced and at the time of the ARPANET trenefor to DCA, it was generally necessary to obtain detailed information about the resources at a particular site directly from the site.

Many different kinds of resource sharing use were anticipated at the ingeption of the ARPANET pregram. Of these, the remote use of one pregram by enother was one of the more distant goals. While some of this kind of resource sharing has, indeed, taken place, such use has grown more slowly than other network uses.

On a less agientific and more programmetic note, the initial program plan anticipated a technology transfer of the network to a common cerrier and the evallability of the network to other semmunities of interest such as the NSF and the NIH. The transfer that did eventually take place was to the Defence Communications Agency and although some utilization has been made of the ARPANET by other agencies like the NIH, it proved rather inconvenient to encourage major utilization outside the Defence Department.

)

4. APPLICATIONS AND CONSIDERATIONS FOR THE FUTURE

4.1 Conclusions of Technical Fearibility

The ARPANET project proved the technical feesibility of schieving reliable high performance, cost effective digital communications by means of packet switching technology and, in turn, the technical feesibility of operating a resource sharing computer network based on this technology. The ARPANET project elso proved the feesibility of aghieving closely knit communities of technical interest over a wide spread geographic area; it is possible that this social feesibility demonstration is as important as the meny technical feesibility demonstrations.

4.2 Recommendations on Additional R&D Requirements and Opportunities

The ARPANET program has represented a significant merriage of computers and communications and since it was the first such marriage of this magnitude, it is not surprising that the program has many different kinds of research progeny. These new requirements and opportunities fall in several different classess

- o opportunities for research in the specific packet switching technologies inherent in the ARPANET;
- o opportunities for research conserving other kinds of packet switching networks employing related technologies;
- o apportunities for research in connecting and integrating different kinds of packet switching networks:
 - o apportunities for research in new uses of networks, and new techniques in Host computers for taking advantage of networks;
 - o opportunities for research in related technologies,

Each of these estemories will be separately discusseds

ARRANET BERRILLE RED

At the point in the ARPANET prejets where ARPA began to consider thensiver of the network to the segio of the Defense Communications Agenct, direct ARPA support of research on the specific packet switching techniques in the ARPANET was

negestarily reduced. However, the network itself continued to grow, both in size and utilization and some technical problems sleerly required additional research and development. Perhaps the most pressing single technical problem of this type was the "routing" problem; and there is currently a given need for improvement in the specific ARPANET routing techniques. Closely related areas like flow control and congestion control would necessarily require remaxemination in gennegion with the detailed review of the routing issue.

Othor Packet Aulteblug Techiologies

In the gourse of research an the ARPANET, two other related kinds of peaket switching technology were initially investigated; packet satellite technology and peaket radio technology. It is clear that these related packet switching technology areas represent major research and development opportunities with strong basic requirements flowing from the needs of the Defense Department. In fact, ARPA is surrently sursuing major RED investigations in both of these areas and significant precress is being made.

Another RED opportunity exists in use of peaket switching techniques on "buses" such as coexis: cables, on power grids; some research in this area has already taken place.

Interiotalna

As ARPA research leads to other networks such as packet radio networks, packet satellite networks, and as other packet switching natworks within the Defense Department, in the U.S. public sector, and in other countries come into being, there is a clear and pressing need to understand how to interconnect such networks in a technically adequate and east effective fashions. Again, clear requirements flow from defense needs and major opportunities for progress are available. ARPA is gurrently supporting a major RSD sativity on the investigations of internetting technology and an experimental attempts at internetting.

Exačeličalalaketrak.Uga

Perhaps not surprisingly the largest area of R&D opportunity
lies in the direction of attempting to use this new packet
switching technology in many different ways to the advantage of

the Defense Department, Some such research and development is already under wey and excellent opportunities exist for cooperative RED programs between ARPA and the verious other components of the Defense Department.

An attractive RED opportunity exists in the investigation of speech transmission over packet switching networks. There is an extremely serious Defense requirement for energeted speech and the confluence of packet switching technology and improved vectoder technology is providing a major RED opportunity for achieving cost effective digital speech transmission and the related possibility for such speech transmission to be encrypted.

The supprisingly suspessful experimental use of network mail has probably opened a new era in sommunisations in the Defense Department. It is expected that such mail systems built on packet switching networks will have an enermous impact in both the Defense Department and the remainder of the public and private sesters. A major research and development apportunity exists to improve the performance and the soft effectiveness of such mail systems and to experimentally deploy them in the Defense Department.

It is diese that pecket switching networks have a very direct amplication to command and control and a significant research opportunity exists in attempting to investigate such command and control applications. ARPA is already proceeding to develop testbeds whereby packet switching technologies and other related technologies can be experimentally employed in dooperation with one or more of the military services. There is a clear requirement for improved command and control and the packet switching technology developed in the ARPANET provides a major opportunity to make progress towards this goal.

At a deep technical level, the proper design of Host operating systems for efficient and cost effective connection to networks is still in the area of significant research and development opportunity. In the ARPANET project, the network connections were "addmens" and it is clearly time to mount the more desired investigation of how best to ascomplish auch connections.

Related leakenlagy

In the course of the ARPANET program other opportunities for research, and development were observed, some glosely related to

packet switching technology and some really quite distinct technologies in their own right.

The deneral tepic of computer operating systems security and dommunications sacurity has long been a subject of necessary and development interest in ARPA. The advent of peaket switching technology has permitted a modified viewpoint to be taken in some of these areas of security research. There is now, for example, a greater consurrence in the security community that minimum technology sections are time and then trusted deep within the guts of encryption systems. Second, packet switching technology permits a different viewpoint to be taken on the techniques for key distribution and there are research and development apportunities specifically in this area currently being pursued by ARPA. Thus, packet switching networks have created additional requirements an accurity technology and at the same time provided additional opportunities.

In the course of the ARPANET development, the Pluribus multiprocessor project was initiated in order to previde a high performence nodel switching unit; However, a multiprocessor technology is, in principle, broadly useful across the spectrum

of data progessing. Thus, the developments in multiprocessor technology that took place within the ARPANET program provide an entiting expertunity for further respects and development in the design and construction of multiprocessors for other Defense computational needs:

The teghniques for remote control of computers in the field developed within the ARPANET project are probably more broadly applicable to the management of computer resources in other areas of the Defense Department. These remote management techniques represent enother opportunity which has grown out of the ARPANET experience.

It was earlier indicated that it was net treatable for the Network Information Center in the ARPANET to keep current oneline detailed descriptions of program resources available within the Host community. This difficulty in turn represents a research and development opportunity for the future. Specifically, research and development is required on how to properly describe computer programs and how to create standards for such descriptions such that it will be easier to create compendia of available resources. In other words, despite the great success.

of the ARPANET, the bests resource sharing problem still represents a fortile area for research and development. Now that the networking technology itself is "given", attempts at description and documentation of Most resources might more easily yield to a research and development effort.

S. PROGRAM IMPACT AND ASSESSMENT OF TECHNOLOGY DEVELOPED

5.1 Service Use of Technology

Seldom has ARPA had greater euceess in convincing its client, the remainder of the Defense Department, to edopt a particular technology. Not only was it possible to transfer the ARPANET itself to the Defense Communications Agency (Code 535), but the Defense Communications Agency has already ambarked on the progurement of its primary backbons major communications system of the future -- AUTODIN II -- based very specifically on the packet switching technology developed in the ARPANET. Even beyond this very major step all three services are actively involved in investigations of packet switching technology to their specific needs.

5.2 Impact on Nan-DoD Programs

In just the very short time since the insertion of the ARPANET this technology has elready resulted in the formation of a new industry; a private sector development of "value added" packet switching networks. Two corporations, Telenet and Tymeshers, are surrently marketing packet switched sommunications

to the General Public as common carriers under the Communications act of 1934. It is difficult to recell any other ARPA program which has had such a direct and immediate exploitation in the commercial sector. In addition, all over the world naw communications systems are being designed and built to take adventage of the packet switching technology demonstrated by the ARPANET project. At least three countries we dreat Smitching France and Canada we have major national PTT spensored packet switching networks either eleedy operating or under development and many other countries are equively pursuing this technology.

5.3 Applications of the Program Results

The most specific result of the ARPANET Program has, of course, been the ARPANET itself; and the ARPANET itself is surrently fully specificant under the management of the Defense Communications Agency and is actively serving thousands of individuals on a delig basis.

5.4 Advance in the State-of-themart

The ARPANET program has represented a firsterank advance in the state-of-the-art of communications and the state-of-the-art

of computer technology. The greatest advance has been in the provision of cost effective, reliable, high performance disital communications, but very significance statement therest advances have also taken place in many other areas, such as topological aptimization, routing, multiprocessor technology, protocols for resource sharing between programs, network mell systems, and remote semputer menagements. Computer and communications technology will have be the same.

6. BIBLIOGRAPHY OF REPORTS

There has been an enormous outpouring of literature about the ARPANET in particular, and about resource sharing computer networks in general. Literally an industry has been formed in the space of 2/3 of a decade and the amount of literature reflects this really unusual metabolism.

The initial seminel papers on the ARPANET were presented in May of 1978 at the AFIPS Spring Joint Computer Conference in Atlantic City and are published in the Proseedings of that conference (AFIPS Conference Proceedings, Vol. 36, AFIPS Press, 218 Summit Avenue, Mentvale, New Jersey 87645), Another early ARPANET session was held at the 1972 Spring Joint Computer Conference in Atlantic City and these papers are published in AFIPS Conference Proceedings, Vol. 48, These two early sets of papers represent a sensible introduction to the early notions and plans for the ARPANET.

A sizeable bibliography and index to publications about the ARPANET was published in 1976 with ARPA supports "Salested Bibliography and Index to Publications About the ARPANET", Becker and Hayes Inc., February 1976, 185 pages, AD=A026980. This

doqument represents the most complete svelishe listing about ARPANET publications, but it is a bibliography and against provide help in trying to selest which of the many references would be suitable to look at in what order. Another bibliography on the literature of resource sharing computer networks in general (not just the ARPANET) has been issued by the U.S. Department of Commerce, National Sureau of Standards: "Annotated Bibliography of the Literature on Resource Sharing Computer Networks", NBS Special Publication 384, revised 1976. This document is also not very useful in helping the reader decide what is important and what is important and what is important and what is important

Severe) velumes of reprints have been published which deel with computer networks in general (rather than the ARPANET specifically), but which attempt to sellest together the important pepers themselves and provide a much easier entry to the literature than the large bibliographies? (1) "Advances in Computer Communications", Wasley W. Chu, Artech House, Inc., 618 Washington St., Dedhem, Massachusetts 82686, 1974; (2) "Computer Communications", edited by Paul E. Green; Jr. and Robert W. Lucky, IEEE Press, 345 East 47th Street, New York, New York 18817, 1975; and (3) "Computer Networking", edited by Robert P.

Bland and Ira W. Cotton, IEEE Press, 345 East 47th Streat, New York, New York 18817, 1976.

A very helty, but fairly readable compandium with a vary large list of references is "Infotesh State of the Art Report 24, Network Systems and Software", Infotesh International Ltd., Nicholson House, Meidenheed, Berkshire, England, 1975. This document deals with far more than just the ARPANET, but it tries to put the ARPANET in context with other surrent work as of 1975.

There are three important reference decuments which have been prepared for the Defense Communications Agency by the Network Information Center at Stanford Research Institute, Menio Park, California

94825 which are specifically addressed to various assects of the ARPANETS (1) "ARPANET Resource Handbook", NIC 39335, Desember 1976; (2) "ARPANET Protect! Handbook", NIC 7184, Revision 1, April 1976; (3) "ARPANET Directory", NIC 36437, July 1976;

A tembook on computer networks (and almost the only useful one so far) is Davies and Barber, "Communication Networks for Computers", John Wiley, 1973.

There have been several regent conference and conference proceedings dealing with computer networks and some selected recent papers from these conferences concerning the ARPANET are as follows:

...

7. DISTRIBUTION LIST

The ARPANET Completion Report should be made publically available through NITS, and copies directly sent to the followings

Former Directors of ARPA involved in the ARPANET efforts

Chaples M. Hersfeld Eberhardt Rechtin Stephen J. Lukasik

Former Directors of ARPANET/IPTO involved in the ARPANET efforts

J.C.R. Licklider Iven Sutherland Robert Teyler Laurence Roberts

Former ARPANET Program Managers:

Stephen Crosker Craid Fields

Some individuals with special interests

Paul Baran Thomas Mariii Donald Davies Leonard Kieinrock Heward Frank Douglas Engelbert Frank Heart

Current ARPANET Spensors Group!

Naval Ship Research and Development Center Defense Communications Agency we ARPANET Management Branch NASA Meadquarters

Pormal Completion Report

U\$ACC aupport Branch
National Sequrity Agency
Defense Communications Engineering Center
Air Force Systems Command Meadquerters
Energy Research and Development Administration
National Bureau of Standards
Advanced Research Projects Agency

Selected DoD agencies and offices:
Other selected government agencies:
Selected congressional committees:
Selected civilian institutions:

CHAPTER III: A REVIEW OF THE ARPANET PROJECT

1. HISTORY

The ARPA Computer Network, or ARPANET as it has some to be delied, consists of IMPs, lines, and hosts as shown in Figure 1. The IMPs (short for Interface Message Processors) are small, special-purpose computers connected to sach other by telephone lines. The hosts are a heterodeneous sellection of computers

Figure 1 == IMPs, Lines, and Hosts

used for a veriety of applications. The IMPs provide a dommunications subnetwork through which hosts communicate with such other, much as the commercial telephone system prevides a communications subnetwork through which humans communicates. In other words, the IMPs and lines make up a communications utility of which the hosts are users. Each host is connected to one IMP; one IMP may have several hosts connected to it.

When deta is to be sent from one host to enother, the date is broken into discrete entities called "measage" at the sending host, which is also estiad the "source host". Each message consists of some data and an address. The address specifies to which host the data is to be sent, that is, the "destination host". Successive measages are passed from the source host to its IMP, where they are broken into smaller entities gailed "packets", each of which carries the same address as the message. The paskets are routed from IMP to IMP screen the network until they arrive at the IMP to which the destination host is connected, where the original messages are reconstructed from the paskets, and passed to the destination host.

The ARPANET was senseived in the middle to late 1948s as a project to be spensored by the Information Processing Techniques

DRAFT History

Office of the Advanced Research Projects Agency (ARPA) of the U.S. Department of Defense (DeD). Construction of the network began in 1969. By 1975 the network had gone through several stages of development and for all intents and purposes had become an operational Department of Defense computer network. In 1975, control of the network was transferred from ARPA to the U.S. Defense Communications Agency, on agency better suited to the administration of a working facility.

The ARPANET is a major development in the evolution of computer communications. Our purpose here is to present the history of the ARPANETS why it was built, who helped build it, the major design decisions (and sems of the minor ones) associated with it, its evolutionary development, its meturity, and why it is important.

1.1 Background

Before relating the history of the ARPANET itself, it is appropriate to summerize the events and elegymeteness out of which the ARPANET same into being. What is ARPA? What was the state of computer communications technology prior to the ARPANET? From what other projects is the ARPANET descended? These questions we enswer in the following subsections.

1,1,1 Capsule History of ARPA and Its Information Processing Techniques Office

ARPA was formed as an agency under the Secretary of Defense as part of the U.S. reaction to Russia's lounch of Sputnik in the midel950s. Soon after ARPA's formation, the U.S. space program was moved out of the military to become NASA, and ARPA focused on the areas of ballistic missiles, solid missile prepaliants, and meterials acience. From these areas ARPA gradually expanded its interests to a veriety of senetimes eseteris fields in pursuit of the kind of ferenceching research projects that gould truly be dailed advanced. The directors of ARPA have been Roy Johnson (1958a1959), General Augtin Setts (1968), Dr. J.P. Ruine (1961a1963), Dr. Robert L. Sproull (1963a1965), Dr. Cherles M. Herzfeld (1965a1967), Dr. Eberhardt Rechtin (1967a1978), Dr. Stephen J. Lukasik (1971a1978), and Dr. George H. Heilmeier (1975apresent). 4a amtence characterizing each directors.

In 1961, Commend and Centrol Research (CCR) was assigned to ARPA. An initial task of CCR was to espitalize on the government's investment in the Gall computers; software, and people at the Systems Development Corporation (SDC), Dr. Ruins chose Dr. J.C.R. Licklider to head the CCR effort and Licklider

loined ARPA in October 1962. In taking the job, Licklider saw that improvements in sommand and sontre! would be heavily dependent on fundamental advances in computer technology, and he was committed to seeking advances in that ifeid, particularly in the subfield of interestive computing. Lightider changed the focus of CCR from a prior relience on SDC to the setting up of research contracts with the best ecademic computer centers, including SDC. Prophetically, Licklider nicknamed the arous of somputer specialists he gethered tegether the "Intergalactic Network". By the beginning of 1964 the CCR program had developed into a farereaching basic research program in advanged computer technology, and Licklider's office was renamed information Processing Techniques (IPT), following Licklider (1962=1964) as Director of IPT have been Dr. Iven Butherland (1997=1977), Dr. Robert Taylor (1977=1969), Dr. Lewrence Roberts (1969=1973). Dr. J.C.R. Lightider equin (1977-1977), and Col. Devid Russell (1977-00-sent).

The centerpiese of the new IPT program was Project MAC at MIT, and the focus of Project MAC was its work in time-sharing, in particular the development of the first time-sharing system depairs of supporting many users. This system was known as the Compatible Time-Sharing System or CTSS. By 1964 CTSS had proven

such a success that mians were being laid for a second generation system which became known as MULTICS. By the late 1962's and serly 1972's, largely because of ARPA support of timescharing research and development of several early timescharing systems, time-charing had become an accepted component of the computer industry.

Time-sharing was only the most prominent of accemplishments beginning to derive from IPT programs during the 1963-1965 period. There were important developments in gemputer graphics, such as the RAND Tablet, Important contributions to programming techniques and languages (e.g., LISP) were made, Support was provided for research in display techniques, computer architecture, and servicial intelligences, by 1965 ARPA had attained a reputation as the deminant source of support for truly advanced information programing research.

In the 1965-1966 period IPT began to shift its emphasis from research to development, By 1967 MULTICS was under development, ARPA had committed itself to the development of ILLIAC IV, and the investigation had begun into means for intersennecting geographically separated computers. Numerous other research efforts were also supported, especially in the software area.

The 1967=1970 period saw the ARPANET effort firmly underways. From that period to the present ARPA has dentinued to support a number of advanced efforts in demputer research and developments.

Over the years ARPA in general and IPT in particular have achieved a number of notable successes. They have had a knack for sensing when the time was ripe for a perticular development and then againg quickly to undertake that development; and they have been able to support developments with sufficient money, IPT has probably centributed more than any other semanter science laboratory in the country (one is tempted to say "then all other laboratories") in most advenced areas of information processing.

Section 1

DRAFT History

1.1.2 The State of the Art of Computer Communications
Circa 1967

«this section will be written for the second draft»

1,1,3 The RAND Study of Distributed Communications Networks,

One of the most important warly studies of computer networks was performed by Paul Baran and his salleagues at the RAND Corporation in the marly 1960s, Many condepts central to the leter development of the ARPANET and other demputer networks were First described in the series of reports published by RAND in 1964 (a list of these reports it given in the bibliography at the and of this subsection). These ideas include the improved reliability of a distributed network structure over a sentralized or star network and over somealled decentralized metworks made up of a gollestion of smaller star networks. Extensive studies were undertaken, including simulation of some orid networks. determine how "survivable" a distributed network sould. 98 expected to be after heavy node and link failupes. This study was particularly conserned with the question of keeping a high percentage of the network available and performing well in the face of enemy attacks on the network, from the point of view of its suitability for Department of Defense applications,

In apequifying closely the engineering details of what was salled the "Distributed Adaptive Message Blook Network", Baren anticipated many of the developments in practical networks that

same a full decade later. In the Distributed Adaptive Message Block Network, a "multiplexing station" connects up to 1824 terminals of widely differing characteristics. Automatic useretoeuser styptography is integrated into the network switching technique to ensure efficiency. Both satellite links and low-cost migroways relay systems are suggested as techniques for providing the network with very high data rate singuits. The concept of a "message block" is introduced; a packet of up to 1824 bits of header and data, which is the unit of data transferred in the network. One of the most interesting expects of this study is that it consided that a large-seein digital transmission network was not only feesible but also eastwoffestive, and propaged that many of the switching functions be implemented in hardware. Baran was considering ways of making extremely reliable networks, and se preferred simple solutions and reliable hardware where possible.

The fellowing bibliography includes the entire set of eleven reports in the original RAND study as well as two published papers resulting from that set and from two later reports. An annotated version of this bibliography is insluded in "Adaptive Routing Algorithms for Distributed Computer Networks" (John M. McGuillan, BBN Report No. 2831, May 1974, pp. 14-17).

Bibliography

- P. Baren, *On Distributed Communications: Is Introduction to Distributed Communications Networks, * RAND Corp. Memo RM+3429=PR, August 1964, 37p.
- 8, Books and P. Beren, "On Distributed Communications: II, Digital Simulation of HetePotate Routing in a Broadband Distributed Communications Network," RAND Corp. Memo RM=3103=PR, August 1964, 49p.
- J_s W_s Sm(th, "On Distributed Communications: III, Determination of PathaLandths in a Distributed Network," RAND Corp. Homo RM=3578#PR, August 1964, 71p.
- P. Baran, *On Distributed Communications: IV. Priority, Precedence, and Overload, *RAND Cosp. Memo RM=3636=FR. August 1964, 63p.
- P. Seren. "On Distributed Communications: V. History, Alternative Approaches, and Comparisons," RAND Corp. Memo. RMm3897mPR. August 1964, 51p.
- P. Baran, #Gh Distributed Communications: VI. Mini-Cost Microwave, # RAND Corp. Memo RM=3762=PR, August 1964, 101p.
- P. Garen, "On Distributed Communications: VII. Tentative Engineering Specifications and Preliminary Design for a High-Data-Rate Distributed Network Switching Node," RAND Corp. Memo RM-3763-PR, August 1964, 85p.
- P. Baran, "On Distributed Communications: VIII, The Multiplexing Station," RAND Corp. Memo RMe3764mPR, August 1964, 103 p.
- P. Beren, "On Distributed Communications! IX, Security, Secrety, and Temperofree Congiderations," RAND Corp. Home RMs3765-PR, August 1964, 39p.
- P. Baren, "On Distributed Communications: X. Cost Analysis," RAND Corp. Memo RM=3766#PR, August 1964, 21p.
- Pa Baran, "On Distributed Communications: XI, Summery Overview," RAND Corps Memo RM-3767-PR, August 1964, 23ps

- P. Baren, "On Distributed Communication Networks," IEEE Transactions on Communications Systems Vol. CS-12, Merch 1964, pp. 1-9.
- B. W. Boshm and R. L. Mebley, "Adaptive Routing Techniques for Distributed Communications Systems," RAND Corp. Memo. RM=4781=PR, February 1966, 78p.
- B. W. Boehm and R. L. Mobley, "A Computer Simulation of Adaptive Routing Techniques for Distributed Communications Systems," RAND Corp. Memo RM-4782-PR, February 1966, 35p.
- B_s W_s Boenm and R_s is Mebley, "Adeptive Routing Techniques for Distributed Communications Systems," IEEE Transactions on Communication Technology, Vol. COM=17, No. 3, June 1969, pp. 348=349.

Section 1

1.1.4 The Lincoln/SDC Experiment

In 1965 Thomas Marill and his adileaques at Computer Corporation of America (in Cambridge, Massachusetts) were commissioned by M.I.T. Lincoln Laboratory to atudy the doncept of computer networking. The primary technical sontest at Lincoln Laboratory was Lewrence Roberts, who played an aggive role in the network studys and the contract was, in feet, a subsontract under the Laboratory's ARPA contrest. The study was done in late 1965 and a report on the study was issued in 1966 ("A Cooperative Network of Time-Sharing Computers, Thomas Marilly Computer Corporation of America, Technical Report No. 11, June 1, 1966; a later paper of the same name authored by Thomas Marill and Lewrence Roberts elso appeared in the Presendings of the AFIPS 1966 Spring Joint Computer Conference, pp. 425-451). The report exemined the besit idea of computer networking, considered the available communications techniques and software problems, and recommended that a threewsomputer experimental network be The report suggested linking three existing constructed. computers, the AN/FSQ=32 at Systems Development Corporation, the 18M 7094 et MIT's Project MAC, and Lincoln Laboratory's TX=2. Later in 1966, CCA received enother contract to carry out the finking of the Q=32 and the TX=2. The Q=32 and TX=2 were in fact linked together, and the link was demonstrated. Later a small Digital Electronics Corporation machine at ARPA was added to this network, by now known as "The Experimental Network". It is noteworthy that The Experimental Network linked host computers directly, and did not use IMPs.

1,1,5 The NPL Data Network

Another early major network development was undertaken et the National Physical Laboratory in Middlesex, England, under the loadership of D. W. Davies. The broad system design of the NPL Data Network, as it was called, was first published in 1967, and bears a resemblence to the network proposed by Peul Baran at RAND, and to the ARPANET. The NPL Data Network was specified to be a packeteswitching network and was to have a hierarchical structure. It was proposed that "local networks" be constructed with "interface computers" which had responsibility for multiplexing among a number of user systems and for communicating with "high level network". The latter would be constructed with "ewitahing nodes" connected together with megabit retections in the systems and some connected describer with megabit retections.

While there was considerable technical interchange between the NPL group and those who designed and implemented the ARPANET, the NPL Date Network affort appears to have had little fundamental impact on the design of the ARPANET. Such major expects of the NPL Date Network design as the standard network interface, the routing elgerithm, and the software structure of the switching node were largely ignored by the ARPANET designers.

There is no doubt, however, that in many less fundamental ways the NPL Date Network had effect on the dealen and evolution of the ARPANET.

Considerable detail about the NPL Data Network may be found in the textbook, written by two of the network's designers, entitled <u>Consuntations</u>, <u>Networks for Consultations</u>, <u>Networks for Consultation</u> (D.W. Davies and D.L.A. Barber, John Wiley & Sons publishers, 1973). A bibliography of published papers resulting from the NPL Data Network effort follows.

8(b) teerephy

- D, W, Devies, K, A, Bartistt, R, A, Scentisbury, and P, T, Wilkinson, MA Digital Communication Natwork for Computers Giving Repld Response at Remote Terminels, ACM Symposium on Operating Systems Problems, October 1967, 176.
- K. A. Bertlett, R. A. Scantlebury, and F. T. Wilkinson, MA Data Communication Network for Registine Computers, W. IEEE International Conference on Communications, U.S.A., June 1968.
- D. W. Davies, "Communication Networks to Serve Repid Response Computers," Proceedings of the IFIP Congress 1968, Vol. 2 Hardware Applications, pp. 658-658.
- K. A. Bertiett, "Trenemissien Control in a Logal Date Network," Proceedings of the IFIP Congress 1968, Vol. 2 Hardware Applications, pp. 784-786.
- D. W. Davies. "The Principles of a Date Communications Network for Computers and Remote Peripherals," Proceedings of the IPIP Congress 1968, Vel. 2 Herdwere Applications, pp. 789-715.

- R. A. Scantlebury, P. T. Wilkinson, and K. A. Sartlett, "The Design of a Message Switching Centre for a Digital Communication Network," Proceedings of the IFIP Congress 1968, Vol. 2 = Hardware Applications, pp. 723-733.
- P. T. Wilkingon and R. A. Scantlebury, "The Control Functions in a Local Data Network," Proceedings of the IFIP Congress 1968, Vol. 2 Handware Applications, pp. 734-738,
- D. W. Dayies. "A Versatile Data Communication Network," Collegue International sur la Teleinformatique, Paris, March 1969, pp. 461-467.
- D. L. A. Barber, D. V. Binke, and R. A. Scantinbury, "Implementation of the British Standard Interface for Data Exchange between Sources and Acceptors of Digital Data," Colloque International sur la Teleinformatique, Paris, March 1969, pp. 807-816.
- K. A. Bertlett, R. A. Scentlebury, and P. T. Wilkinson, "A Note on Reliable Full-Duplex Transmission ever Helfeburiex Lines," CACM, Vol. 12, No. 5, May 1969, p. 268,
- D. L. A. Barber, "Experience with the Use of the British Standard Interface in Computer Peripherals and Communication Systems," ACM Symposium on Data Communication, Pine Mountain, October 1969,
- R. A. Seentlebury, "A Model for the Legel Area of a Date Communication Network Objectives and Hendwerm Organization," ACH Symposium on Data Communication, Pine Mountain, October 1969.
- P. T. Witkinson, " Model for the Legal Area of a Date Communication Network a Saftware Organisation," ACM Symposium on Data Communication, Pina Mountain, October 1969.
- D. L. A. Barber and D. W. Devies, "The NPL Date Network," Proceedings of the Conference on Laboratory Automation, Novosibinsk, October 1978.
- D. W. Device. "The Control of Congestion in Packet Switching Networks," Proceedings of the Second ACM/IEEE Symposium on Problems in the Optimization of Data Communication Systems, October 1971, pp. 46-49; also in Transactions of IEEE, Vol. COM-20, No. 3, June 1972.

- R, A. Scantiebury and P. T. Wilkinson. "The Design of a Switching System to Allow Remote Access to Computer Services by Other Computers and Termine! Devices," Proceedings of the Second ACM/IEEE Symposium on Problems in the Optimization of Data Communication Systems, October 1971, pp. 1684167; also in Transactions of IEEE, Vol. COMMIN, No. 3; June 1972,
- W. L. Price, "Simulation of Packet-Switching Networks Centralied on Igarithmic Principles," ACM Third Date Communications Symposium, Floride, November 1975, pp. 44-49.
- R. A. Scantlebury and P. T. Wilkinson, "The National Physical Laboratory Date Communication Network," Proceedings of the ICCC Conference, Steckholm, August 1978, pp. 2274238.
- D. W. Davies, "Packet Switching, Message Switching and Future Data Communication Networks," IFIP Congress, Stockholm, August 1974, pp. 147e158.
- W. L. Price, "Simulation Studies of an Iserithmically Controlled Store and Farward Date Communication Network, Proceedings of the IFIP Congress, Steakhelm, August 1974, pp. 151-154.

1.1.6 The Events Loading Directly to the ARPANET

The ARPA pregram pien under which ARPANET development was undertaken was prepared by Lawrence Roberts and approved by Dr. Rechtin, Director of ARPA, during Robert Taylor's tenure as Director of IPT. There is general agreement that Roberts is the individual with the most valid dlaim to being the "father" of the ARPANET. In a note dated July 7, 1977, Robert Taylor recalled the events leading to the undertaking of the ARPANET development up to the time when Lawrence Roberts leined ARPA.

The ARPANET, like most good ideas, was there at least in tenueus form from the beginning, waiting for the right time to be invoked, in 1962-65 Lieklider with tanquestreshaek used characteristic ... the "intermalactic network" to outline a plan for soupling prodemic computing greenties to the SDC Gm32 in order to eid in giving the SDC effort (which Liek inherited from DDRSE) the benefit of some of the experience gained by the universities. Link was smong the first to perseive the spirit of community crosted among the users of the first time-sharing systems. Mis successor, Iven Susherland, responded in 1965 to a UCLA request involving small seale networking research at one site by offering ARPA support to the erestion of a three-heat compus network, but the internetine behavior of the three temputer center directors at UCLA made the situation hopeless and the effer was withdrawn. Early in 1966, after a year in ARPA, it meamed to me that Lick and Ivan had paved the way for the next step. In painting out the community phonomene created, in part, by the sharing of rescurees in one timesharing system, Lick made it easy to think about interconnecting the sommunities. The interconnection of interestive, one inc communities of people was also a dresm of Deug Engelbart whose work I had supported from NASA in the early 60°s and was continuing to support from ARPA.

From Ivan and the UCLA experiences I had to learn once again that the teuchest problems were not the technical ones. If I dould not get some ARPA funded pertidipants involved in a commitment to a purpose higher than "Who is going to steel the next ten per cent of my memory dyales?", there would be no network. By June of 1964 Ivan had moved to Messachusetts to work at Lincoln Leb in the summer and then teach at Hervard. Ivan and I had anseuraged Larry Roberts, then at Lincoln Leb, to pursue a fimited networking experiment involving Lincoln's These and SDC's Gosz. With Ivan's leaving, IPTO was left with myself, a servetory, and a 320 million/year budget. The Director of ARPA then was Charles Hersfeld whe, sithough a physicist, was extremely imaginative, approachable, and like his predecesor, Robert Sproull, siways interested in discussing new ideas. (It was a pleasure to work for them!) I went to Charles for help. I told him I wented to initiate a networking project and I outlined its objectives.

- 1. To enable users at one IPTO spansored site to gain interactive assess to the pregrams/date at eather IPTO spansored sites.
- 2. To face deliberately the problems of intersenmenting incompetible systems in order to someday enable the DOD to eliminate this argument for single-source systems.
- 3. To further enseurage, through suggestful networking, the growth of epselelized areas of semputing knowledge at different epensered sites = a trend which was already developing = thus enabling a uper to tap verious sources and hopefully reduce future duplication of software development.

In addition I gave Charlie come simple, hoped for, performence specifications, sage, no user should have to wait more than one second for prosendountry respense else the illusion of being a local user, near the computer, be destroyed, (I underestimated what I new believe is one of the most important results from the ARPANET experimental the power, and utility of message systems.) Charlie seid "Great! Do it!" and gave the project more than ample budget to get off the ground, I told him I needed some help

and described Larry Reberts as an ideal manager for the project. Charlie approved that request as well, ...

Next I had to denvinge Larry to head the project and convince some currently funded contractors to perticipate 48 users. It may seem surprising but I think the former was the most difficult. I had tried to interest Larry in ARPA estilet in 1966, when Iven had ennounced his plans to leave, but Larry was non-gommittal. After getting Herzfeld's approved for the project I visited Larry several times over the next few months trying to convey the extent to which I needed help, but he argued (as meny have whem ARPA has asked for help) that he did not went to make a break in his research career at that time, furthermore, he reminded me that he was helping support IPTO's coals at Lincoln, which was gertainly true. I returned to Washington ofter one of those visits with Larry and was about to give up. I began building a list of other candidates. Then one day something reminded me that ARPA supported 51% of Lincoln Laboratory® budget, I went to see Charite, told him I was having difficulty getting help, and asked him if he would gall the Director of Lincoln Lab and suggest to him that it would be in Lingoln Lab's best interest if they sould spare Dr. Roberts for an important ARPA project. Charlie placed the call in my presence, Larry accepted the jobs soon after, and moved to Virginia during the Christman holidays of 1964. That assemplishment was probably the hardest and the most important of my ARPANET involvement.

Since 1963, a recurring question about some prepased or existing ARPA/IPTO project put to IPTO Directors by ARPA Management is "Why den't we rely on the computer industry to do thet?"; or accessmelly more strangly, "We should not support that effort because ASC (reed, "gemputer industry") will do it a if it's worth deing;" On the other hand, ARPA's support of research on networking, timesharing, and interactive computing in general, apparently denies the

^{*} Roberts first went to ARPA as special assistant to the Deputy Director of ARPA. His responsibility was to provide technical advice to both IPT and other ARPA offices on computer matters. He later succeeded Taylor as director of the Information Processing Techniques office.

premise. So we ask ourselves if there is some singular, centrel theme which helps us distinguish the ARPA smphasis vis a vis computing research from that of the computer industry.

i believe the ensure is a clear, unembivalent "yes," beautifully undersored by the explicit summunications objectives of the ARPANET. The ARPA/IPTO theme was set by Licklider in his first tenure ("e2e*44) and sentinues to be supported todays it is that the premise offered by the computer (read, "seftwere, hardware, et sla") as a communication medium between people, dwarfs into relative insignificance the historical beginnings of the somputer as an arithmetic engine. The computer industry, in the main, still thinks of the computer as an arithmetic engine. Their haritage is reflected even in surrent designs of their acommunication systems." They have an economic and paychelogical commitment to the arithmetic engine medel, and it can die only clowly; furthermore, it is a view that is atill reinforced by meet of the notion's computer asience progrems.

Even universities, or at least parts of them, are held in the grasp of the erithmetic engine sendent, as an while it may be later then we think a it is certainly going to take lenger than we think; or put more genetruatively, the challenge taken on by IPTC roughly fifteen years ago has unveiled a problem-demain (apportunity-demain) which can keep IPTC productive for seally another fifteen years, even under the ARPA griterion that if it doesn't offer an order of magnitude improvement it should not be supported by ARPA.

1.2 The Events of 1967 and 1968

With Roberts on board at ARPA by the end of 1966, the ARPANET development effort got under way in sernest. The years 1967 and 1968 were spent promoting interest in the ARPANET project both within the government and with IPT sentrestors, deciding the fundamental structure of the network, writing a request for quotation, selecting a contractor, and other related activities.

Each year until recently ARPA held a meeting of the Horingipal investigators from each of its university and other contractors. Called the HPI meetings, these meetings provided a forum at which contractors described their past year's results for the benefit of the other sentrectors, the IPT staff called the contractor's attention to problems of DoD soncern, and everyone together attempted to unsarth areas of technology which needed only ARPA fertilization in order to risen to a seint of usefulness. One of the topics of the PI meeting held at the University of Michigan in early 1967 was networks. At the meeting it was agreed that work sould begin on the senventions to be used for exchanging messages between any pair of computers in the proposed metwork, and also on consideration of the kinds of

dommunications lines and data sets to be used. In particular, it was decided that the intershoot gommunication "protogo!" would include conventions for sharester and block transmission, error checking and retransmission, and computer and user identification. Frank Westervelt, then of the University of Michigan, was picked to write a position paper on these areas of communication, an adapta "Communication Group" was selected from among the institutions represented, and a masting of the group acheduled.

The pien considered throughout much of the Michigan meeting was to connect all of the computers by phone lines and data sets so as to allow any somputer to actabilish communication with any other computer using a linear witching technique (i.e., calling it up on the telephone). A small program in each computer would interface to the data set and phone line and when this small program was given a message to enother computer, the small interface program would perform a "message-switching" and transmission function, deciding how to actually reach the other computer and transmitting the message to it. Hes Clark, then of Washington University, is credited with proposing at some point in the meeting that a small computer be inserted between each participants computer and the phone line. This gangept was

further discussed with the ARPA staff and one or two others (n a now fabled tax) Fide to the airport after the meeting, and the concept of the Interface Massage Processor or IMP was reported in a mamo from IPT to the PIs summerizing the Mighigan meeting.

The IMP would perform the functions of dialaup, error checking, retransmission, routing, and verifications on behalf of the participents? computers (which we shall hereafter call hosts). Thus the IMPs plus the telephone lines and data sets would constitute a "message=switching network". The protocols which were to be established would define the semmunications formets between the IMPs. The interfece between a host and an IMP would be a digital interface of a much simpler sort requiring host consideration of error checking, retransmission, and It was olderly noted that the major disadvantage of routing. inserting the IMP was the dost of installation of another computer beside each host. The major advantage of (negrting the IMP bes inspess be the fact that WAR te a untited. straightforward negwerk design eauld be made and implemented without undue consideration of the varietions and constraints of the host computers. Further, as the network evolved, it would be much simpler to modify the network of IMPs them to modify all the host computers. Finally, the IMPs would relieve the house of the dommunications burdens they were initially scheduled to carry, It was also noticed that if necessary IMPs could be located at strategic connection points within the network to concentrate messages over crossessuntry phone lines, a network of IMPs was likely to be implemented feater than a network of directly connected hosts, and the network of IMPs provided a distinct network entity which would be useful in presenting the network publicly.

The initial draft of data communications protocols was drawn and circulated by Prank Mestervals and subjected to intensive review by a number of interested particle at a meeting at ARPA in midwMay, 1967, Policying that meeting further work on the conventions was carried out by several individuals from various institutions and another draft, by Mestervals and Milis, was directled after the summer of 1967, On the physical side this draft specified full duplex binary serial transmission at a minimum rate of 2488 bits per second in each direction, and further specified that each site should be capable of automatic anaworing of incoming sells and automatic dialing to originate outgoing calls. On the logical side the draft specified ASCII, use of 16-bit CRC sheeksums, a certain communications syntax, use of either a text or binary formet (the letter using DLE-doubling

for transparency), and sertain centro; conventions. This draft also noted that the concept of using IMPs was gaining support.

A two-day meeting was acheduled in early October 1947, at ARPA, to discuss the protocol paper and specifications for the IMP (which had experently become the accepted method). The meeting anneuncement else included a survey to be filled out by each of nineteen potential network locations providing cross quesses of total ARPANET traffic to other network locations by mid-1969 and percentage of the total traffic to each other location. Sy this time the presented network was known as the ARPA Network, rather than the vague descriptions that had been used previously.

A veriety of tepics were discussed at the Datober meeting including message formating, message protocols, dynamic routing and message propagation, queuing, error control, measurements, and IMP-to-host communication. In a memo Reberts noted that the discussion at the meeting helped him formulate many plans and make several decisions. It was decided that 50 Kb communications lines would be used because of the vestly improved response time which could be obtained with these lines as esposed to the proviously proposed 2.4 Kb lines. The 55 Kb lines were to be

leased, eliminating the slow disleup prosedure. The neture of the telephone teriffs available to the government made use of 58 Kb lines compensively inexpensive. With each IMP normally permanently connected to three other IMPs via the leased 50 Kb lines, the IMPs were to use steremend-loowerd techniques to provide fest message handling, Each IMP was to accept of up to 5,000 bits from its host computer and to break this inte 1,988-bit peckets. Each pecket was to be treated independently and routed on one of the three intermIMP lines. When a packet arrived at an IMP, it was to be stored, arrop sheeked, and routed on to a further IMP. At its destination, packets would be held until an entire message sould be assembled and then delivered to the destination host. With three lines per IMP, it was expected that approximately three storemendsforward stages would necessary to get a message from any one of twenty locations to any other. It was also decided that a lambit CRC would be computed for each message, possibly by a special hardware attachment to the IMP. Also, an aight-bit binory format was to be used, requiring the IMP progrem or special hardware to insert Diffecharacters before control characters and strip them out: Finally, it was suggested that the dennection between the IMP and host use a single seriel line between a shift register in the IMP and a shift register in the host, with the shift registers matched to the word sizes of the IMP and host. This series sine was to be expedie of operating at 8 Mb. The need for control lines between the IMP and host was also noted.

The first half of 1968 was a time of intense study by individuals of the methods to be used in the ARPANET. example, in Jenuery 1968, Dr. Robert Kahn wrete a meme to "ARPA Network Committee" in which he ergued that the mean time between undetected errors in the network should be at least an order of magnitude larger than the debugging time for the network, and he condiuded that an error detection scheme using a B.C.H. sade generated by a 24-bit shift register be used. In February, Leonard Kielnrock of UCLA elregisted three working papers about the ARPANET: a note by John Stehura on a suggested network control language which Kleinrock characterized as a first effort at salving a critical problem from the users, point of views a note by Kielnrock presenting his thoughts on a simulation of the ARPANET; and a preliminary quauling enalysis by Kieinrock for storemendeferward response time of the ARPANET. Other UCLAugenerated documents on the ARPANET followed. In early March Kahn distributed another memo in which he urged that an effort be made to obtain data on the nature of errors on the 50 Kb communications lines, suspected that an important objective of the natworking experiment be the evaluation of the intersection between communication feellities and terminal data processing equipment, and suspected that the initial network include at least one transcentinensal link.

In perallel with the individual research that was going on, ARPA moved to pull together an IMP specification which could serve as the basis for a work statement for a sompotitive procurement to select a contractor to build the IMP network. At the end of 1967 ARPA initiated a small centract with the Stanford Research Inactivate for the development of specifications for the necessery communications system. Eimer Shapiro was to be the key person on this study. Published in a final version in December of 1968 as a Timpeds SRI Report entitled "A Study of Computer Network Design Peremeters, on early version of this report in early 1965 served on the first draft of the IMP specification. (Incidentally, this study was the first bit of explicitly paid affort doing towards the ARPANET; all effort prior to that had been previded by contractors working extra hours on their existing ARPA sentrests.) In February or Hargh & meme written by Shapiro and revised by Kleinreck entitled "Functional Description of the IMPR was direulated. After the first dreft by Shapiro, it

is believed that Glann Guller wrote a second draft, and Reberts and Wessier of ARPA wrote the final version of the IMP specification. In any case, by the first of March, 1968, IPT was able to report to the Director of ARPA that specifications for the IMP were essentially complete, and that they would be discussed at the upcoming PI meeting with the goal of issuing a Reducat for Quotation shortly thereofter. The network was discussed at the PI meeting and by June 1968, the ARPANET procurement efficially started.

To begin an ARPA program, an ARPA Program Filen is needed, This plan must be signed by the director of the ARPA office involved, in this case IPT, and the director of ARPA, No one else in the DoD had to approve the plan to presure the ARPANET, In fact, it appears there was no serious emposition to the ARPANET within DoD at large, and within ARPA the most negative word was an admenition from the Deputy Director, Dr. Lucesik, to so slow and keep it small. In fact, as the incoming program manager for the ARPANET, Roberts had only to win over Director Rechtin (and his Deputy, Lucesik) and IPT Director Taylor. It has already been noted that Taylor was instrumental in getting Roberts to ARPA in the first place for the purpose of managing a network effort, and with Roberts then serving as Special

Assistant for Information Sciences reporting directly to the Deputy Director in the Director's office, it would have been surprising if the Director had not conquered with Roberts's plans. The ARPA Pregram Pion for the ARPANET, entitled Resource Sharing Computer Networks's, was submitted June 3, 1968, and approved by the Director June 21, 1968. The fixed year 1968 ARPA budget included approximately \$500,000 cormerted for the ARPANET effort.

The program plan for the ARPANET is an interesting document, the stated objectives of the program were to develop experience in interconnecting computers and to improve and increase computer research productivity through resource sharing. Technical needs in adjentific and military environments were cited as lustification for the program objectives. Relevant prior work was described. It was noted that the computer research senters supported or pertially supported by IPT provided a unique testbed for domputer networking experiments, as well as providing immediate benefits to the centers and valuable research results to the military. The network planning that had gone on was described, the need for a network information senter was noted, and the network design was sketched. A five year schedule for network progurement, construction, operation, and transfer out of

ARPA was presented; (It is noteworthy that IPT initially had in mind eventual transfer of the operations) network to a common service;) Finally, a severel-millian-dollar, severel-year budget was stated.

The Defense Supply Service a Washington (DSSaW) agreed to be the procurement agent for ARPA. At the end of July the Request ing Quotation for natwork IMPs was mailed to 148 potential had expressed interest in reselving it. biddess who Approximately 188 paople from \$1 companies attended a subsequent bidders' conferences: Twelve proposals were estually received by D85=W comprising 6.6 edge=feet of paper and presenting an avecome evaluation task for IPT, which more normally awards contracts. On a sole source basis. Attempting to evaluate the proposals Matrictly by the book", an ARPAwappointed evaluation committee retired to Monterey, California, to carry out their task. ARPA was pleasantly supprised that severel of the respondents believed that they sould construct a notwork which performed as much as . B factor of five bottom than the delay constraint given in the RPG.

^{*} In parally! with the competition for the IMP contract, ARPA hald a competition among the common confiders to obtain the communication lines for the initial network. The contract was awarded to ATET by the Defense Commercial Communications Office in early September, 1968,

Four bidders were reted within the sone of sontention to receive the IMP centract, and supplementary technical briefings were requested from each of these bidders. Final negotiations were serviced out with two finalists, and one was sheem in the week before Christmes, 1966. The sontract was swanded and work began the second day of the New Year in 1969.

1.3 Key Aspeats of the RFG

The RFG apecified both administrative and technical aspects of the contractor selection and project implementation efforts. Bidders were requested to quete on a "goot plus fixed fee" (CPFF) basis. This is an arrangement the government summonly uses when contracting for unique developments, wherein the bidders first make their best estimate of the beforesprofit sost of sempleting proposed effort. The bidder and the government then megatiate an appropriate fixed fee taking into eccount the beforesfee estimated sost and the risk the bidder is taking in accepting the contract. If the contractor leter overrung his estimated cost, the government has the estion of stepping the effort before completion or paying the additional cost to demplete: hewever, no edditional fee is peld on the additional costs to complete. A CPFF errangement makes sense when the government wents a ricky development done and would have trouble getting a contractor to undertake the development if the contractor had to been too much of the risk himself, as would be the case with standard fixed price bids. On the other hand, the fixed fee part of the arrangement reduces the risk that the contractor might artificially run up his equat, since he will gain no additional fee.

It was apacified that responses to the RFG would be evaluated on four aritaria in addition to east:

- 1. Understanding and depth of analysis of sechnical problems involved.
- 2. Availability of qualified, experienced personnel for axeignment to aptimere, hardware, and installation of the system.
- 3. Estimated functions; performence and chaise of hardware.
- 4. General quality, responsiveness, and corporate commitment to the network concept.

These four criteria were given weighted values of 38, 25, 25, and 20.

The RPG had provision for a bidders conference and stated that there would be no other opportunity for bidders to dispuss technical issues with the government. Bidders were asked to provide a system design for a ninetaen=IMP network, but to price a four=IMP network. A thirteenementh performance period was requested to include design, construction of a pretotype IMP, and implementation and installation of four operational IMPs. The four IMPs were to be installed nine months after start of the contract with the contractor supporting them in the field for three menths after installation. The contractor was required to take full system responsibility, although subcontracting a portion of the work was a possibility.

The statement of werk included in the RFG was subtitled "Specifications of Interfere Message Processors for the ARPA Computer Network" and specified many ARPANET characteristics which some people later thought represented technical decisions on the pert of contractor. While there was room within the guidelines of the RFG for the centractor to exercise his technical judgment in many matters, in fact ARPA knew what it wanted and the basic network decion should be credited to ARPA and those who helped ARPA in 1967 and 1966. The decument was a Request for Quetation for the implementation of ARPA's system decion, not a Request for Presided for an alternative system decion. The following figure reproduces the Table of Contents of the Statement of Herk on IMP specification section of the RFG; it indicates the extent to which ARPA understood exactly what sent of network it had in mind.

The IMP specification steerly delinested the division of responsibilities emeng heat sites, IMP contractor, and telephone company. Each individual host site was responsible for designing and implementing for its own convenience the hordware and software necessary to attach the heat to the natural, and the hordware and software and software to utilize other house on the network. The telephone company was to be responsible for providing necessary

```
Network Description
I.
            Introduction
            Functione1.Description
                The User Subnet
                The Communication Subnet
            Functional Description of the IMPS
                Breaking of Messages into Packets
                Menagement of Message Buffark
            2.
                Routing of Messages
            3,
                Generation, Analysis and Alteration of
                Fermatted Messages
                Coordination of Activities with Other IMPS
            ٩,
                Coordination of Astivities with its HOST(s)
            61
                Measurement of Network Parameters and
                Funetions
                Detection and Disposition of Faults
            8.
                IMP Software Separation Protection
            The HOSTWIMP Interfores
            The IMPECARRIER Interfaces.
            Notwork Porformance Characteristics
               Messes Delay
Reliability
            1.
            2.
                Notwork Capacity
            3,
                Network Model
            HOST-HOST Characteristics
            IMP=Operator Interface
```

II. Network Contractor Perfermence

III. Elements of System Design

Appendix

- A. ARPA Network Nedes
- B. ARPA Network Topelogy
- C. IMP Delivery Schedule
- D. Input and Output Fagilities for the IMP Operator
- E. ARPA Network Data Rates Bayween Nodes in Kijebite/spe,
- F. Date Communications Conventions
- G. Routing

Figure 2: Table of Contents from RFQ Statement of Work

directs, data sets, and line conditioning equipment utilized by the network. The IMP contractor was to be responsible for providing necessary handware and seftware to connect IMPs to each other using the circuits supplied by the telephone company and to connect IMPs to hosts, as well as providing handware and seftware necessary to implement the procedures which allowed dreation of a network of IMPs depable of forwarding messages from one host to enother.

The functional description of an IMP specified the use of messages not leader than \$192 bits which would be braken into packets of net more than 1924 bits. Messages were limited in size to make them mensageable for the hosts. Shorter packets were used to reduce the probability of transmission error with the attendant necessity for retransmission. It was noted that IMP provision of message and packet buffer space would permit speed conversions to take place, provide queuing space in the face of delays, and permit retransmission in the event of erroncous transmission. A require algorithm was hypothesized which would take into assent the connectivity of the network, IMP and line busyness, and message priority, and use this information to forward a packet to the next IMP on a path to the ultimate destination, periodic updates based on exchange of routing and

information with other IMPs and hosts was also loading hypothesized. An IMP was to topodinate its activities with other IMPs and its hosts and perhaps other special hosts. IMPs were to take messages from a lass! host at the IMP's convenience, but to send messages to a local host at the hestes genventages. The IMPs were to be able at selected times to measure selected network parameters and to theme the mayements of selected messages through the network. The data resulting from these measurements and tracings was to be espable of transmission to a host, and the measurement activity was to be capable initiation and termination by a host or enother IMP. The IMPs were to detest and receier from various IMP, host, and line failures. In particular, it was to be sessible to stop, start, examine, or reload IMPs from selected network hosts. Finally, it was thought that at each IMP after it would be possible for sonsial heatuspecific code to be provided by heat site programmers, and thus it was desirable to protect the rest of the IMP from the pertien that the hest personnel sould access and DPOGPam.

There were two perticularly interesting appears of the hostotest interfece was to be hostotest interfece was to be bedified between them a different one for each host; 2) each

bidder was to genelder the sout of providing interfaces to multiple heats per IMP, although only one heat interface was required; and 5) sufficient program space was to be left to do host-apealific character gode donversion and repeaking of binary massages.

The interface between an IMP and (to telephone lines was required to have herdware to sense theracters, detest central sharecters, detest central sharecters, delicate and shaek the 24-bit CRC, provide a real-time clack with 28 migroseconds resolution, and provide fault and status information. Further, the IMP was to be optimized to handle three lines but be capable of handling six.

Several network perfermence theresteristies were specified. The average message delay for a short message to so from a source IMP to a destination IMP was to be less than one-half second for a fully loaded network. The probability of lost messages and message arrors was to be very low. Interestingly, network separity was considered third in order of importance and was defined to be the maximum bit rate that can be input at every nade and still have the message delay remain less than one-half second; a 20 Kb network capacity was hoped for. A network model was presented.

Host-to-host traffic flows were estimated and it was hypothesized that there would be a trimodal distribution of traffic type (high rate and short length, madium rate and madium langth, and law rate and long langth).

In addition to considering the option of multiple hosts connected to an IMP, the bidder was also eaked to consider the provision of memory protection to facilitate simultaneous IMP operation and charkout of new software, and to consider what additional hardware and software would be necessary for an IMP to provide a terminal concentration capability for its host or for the network (i.e., no host, lust terminals).

DRAFT

1.4 Chronological Development, 1969 to 1975

Within a year efter the eward of the IMP contract, the first IMPs were installed in the field. Hosts were connected to these First IMPs and a series of network messurements was undertaken, Heat software was written, heats began to communicate with each other: more IMPs and hosts were added; and gradually the ARPANET became an operating entity. After six years of development and operation the network was no longer well suited to menagement by an accord with the charter to appnoon edyaneed research and development, and thus network management was transferred to the Defense Communications Agendy. In the following subsessions we describe the happenings of the years 1969 to 1975.

Section 1

DRAFT History

1.4.1 The Groups and the Key People

The ARPANET development was a joint effort of many individuals and institutions, all responsive to ARPAPs directions. Before we describe the ARPANET development further, it is best to list the principal "players" and befeffy describe their cross of responsibility.

1,4,1,1 Menagement and Administration of the Ngtworks ARPA IPT, D\$SeW, RML, and DECCO

Neturally, ARPA IPT played a major role in the development of the network, and Lawrence Roberts was the mest important player of them all. Roberts promoted the network, made gertain design designons, led the evaluation team which selected the IMP contractor and selected other involved empteetors, and was the program manager for the network program.

IPT usually done little daystooday management of its centrasters. Especially with its research sentracts, IPT would not be producing faster results with such management as research must progress at its own pass. IPT has senerally edepted a mode of management which entails finding highly motivated, highly skilled contractors, giving them a task, and allowing them to preced by themselves. Of course, IPT holds periodic reviews, and in some cases the selentiate who are serving as IPT Director or program managers sitch in and themselves perticipate in the sontractors research. If a contractor does not produce good nesults in this environment of low-profile management, the dontractor is dropped. If the contractor does produce good results, IPT continues to fund him. Remember that IPT generally

awards contracts on a sole source basis, which makes it possible to drop or keep contractors. The IPT style of management has also been distated by the exceedingly small staff which has been allowed for them for many years, which has governly constrained the time that could be spont supervising each contractor. Within these constraints, the IPT management style has been reasonably successful.

The ARPANET program was managed by IPT in the same manner as most of their research pregrams, despite, the fast that the dentract for the IMP development resulted from a competitive programment and the network evolved (nto an operational entity, Thus, while IPT set policy for the network, made decisions about who would loin the network, adjudicated squabbles between the several involved contractors, and so forth, IPT did not run the network day, BSN provided daysbyseday operation and maintenance of the network, and BSN and the other contractors involved carried out much of the daysbyseday business of the network among themselves without need for constant IPT supervision.

Initially IPT (test) took care of ordering telephone lines, filling out forms which were ment to DECCO for procurement from

and execution by the various telephone companies. Also, IPT initially did its own technical monitoring of the various contractors associated with the ARPANET; DSS=W was used as the precurement agency for the contractors. Eventually, extempting to rid itself of the routine, tedious aspects of ordering communications lines and providing technical monitoring of contractors, IPT shifted the DSS=W functions, routine dealings with DECCO, and some contractor technical monitoring to the Range Massurements Laboratory at Petrick Air Force Base in Flanday in addition to having a presumement genebility, RML also had a technical support sapability.

There were several other members of the IPT staff who have been prominent in the management of the ARPANET besides. Roberts, In 1968, 1969, and part of 1978 Berry Wessier halped Roberts with the teshnical side of ARPANET activities. Until the shift of responsibility to RML, first Col. Bruce Doin and leter Major Dave Caristrom took sare of ordering phone lines. Eventually Mrs. Stephen Crosker became ARPANET program manager. At about the time Licklider returned for his sesond term as IPT director, Drs. Craig Fields took over as ARPANET program manager and Fields was followed by Mr. Stephen Walker (with a brief interregnum by Mr. William Capison). Walker served as program manager until the

trensfer of the network from ARPA, and he continues today to be the ARPA point of dentest for ARPANET conderns. Of source, both IPT directors after Roberts, Lisklider and Russell, had attend influence on IPT's ARPANET activities. Throughout IPT's involvement with the ARPANET, Mr. Al Blue has been involved in the financial aspects of the network's development and operation (Blue was acting IPT director between Roberts and Lisklider and during that period was deeply involved in other espects of the network's management).

As mentioned above, IPT eften played a strong technical role in the ARPANET, and members of the IPT staff wrote a number of papers describing the ARPANET activity. Several of these papers are listed in the following bibliography.

Bibliography

- T. Merill and L. G. Reberts, "Toward a Cooperative Network of Time-Shared Computers," Proceedings FJCC 1964, p.425-431.
- L. G. Reberts, "Multiple Computer Networks and Intercomputer Communication," ACM Symposium on Operating System Principles, October 1967, 79:
- L. G. Roberts and B. D. Wessler, "Computer Network Development to Achieve Resource Sharing," Prescedings SJCC 1978, p. 543=549.
- L. G. Roberts, "A ferward Leek," Signal, August 1971, p.77*81,
- L. G. Roberts, MARPA Network Implications, * EDUCOM Bulletin, Fall 1671, p.4-6.

- L, G. Roberts, "ARPANET's Current Status, Future Plans," EDUCOM Spring Conference, April 1972, p.7012.
- L. G. Reberts, "Network Retionales A Seyear Reeveluation," Proceedings COMPCON 1973, Pebruary 1973, p. 3-6.
- L. G. Roberts, "Dynamic Allogation of Satellite Capacity through Packet Reservation," Proceedings NCC&E 1975, p.711-716.
- L. G. Roberts, "The ARPA Net," in Computer-Communication Networks, N. Abramson and F. Kuo (Eds.), PrentiqueMail, 1973.

1.4.1.2 The Network Analysis Corporation

Led by Dr. Heward Frank, the Network Analysis Carsonation (NAC) of Glan Cove, Long Island, was put under congress by ARPA to specify the topological design of the ARPANET and to analyse its cost, performance, and reliability characteristics.

In the process of evaluating any of the parameters of a particular network design, such as east, reliability, daisy, or throughput, it is necessary to simulate the flow of traffic through the proposed network. Then, the design may be altered elightly to improve one of these measures. In this procedure, it is important to have a facility for specifying the neutral on which traffic will flow in the network, and the precedure must not be too complex since it must be repeated so often in the iterative design process. NAC has developed same very efficient methods for ingremental changes to a shortestepath routing algorithm as the network topology is changed. Further, they have discovered a feater shortestepath algorithm than was previously evaluable, taking advantage of the low connectivities usually present in most precise; communications metworks.

The following bibliography lists much of the sublished work by NAC related to the ARPANET.

Bibliographys

NAC First Samiannual Tachnical Report for the Preject, "Analysis and Optimization of StoremendoForward Computer Networks," June 1978, 6294

NAC Second Semiannual Technical Report for the Project, "Analysis and Optimization of Storecandeforward Computer Networks, " January 1971, 96p.

NAC Third Semiennus: Technical Report for the Project, "Analysis and Optimization of StorewandeForward Computer Networks," June 1971, 1898.

NAC fourth Semiennual Technical Report for the Project, "Analysis and Optimization of Storemend-Forward Computer Networks," December 1971, 123p.

NAC fifth Semiennus! Technical Report for the Preject, "Analysis and Optimization of Storemendsforward Computer Networks," June 1972, 92p.

NAC Fine: Technical Report for the Project, "Analysis and Optimization of Storewand-Forward Computer Networks," October 1972, 1280.

 H_{\star} Frank, I_{\star} T. Friedh, and W. Chou, "Topological Considerations in the Design of the ARPA Computer Natwork," Proceedings SJCC 1978, p.581-585.

H. Frenk, I. T. Friech, R. Ven Siyke, end W. S. Chou, "Optimal Design of Centralized Computer Networks," Networks, Vel. 1, no. 1, 1971, p.43-58.

H, Frank and H, Cheu, "Routing in Computer Networks," Networks, Vel. 1, No. 2, 1971, p.99-112.

H. Frank, "Computer Network Design," IEEE International Convention Regard, March 1971, p.228-221.

^{*}I have asked Howle Frenk to suggest additions and deletions from this bibliography == DCH.

- H. Frank and W. Chou, "Throughput in Computers Communications Networks," International State of the Art Report no. 6; Computer Networks, Infotesh Information Ltd. Maidenhead, Barkshire, England, 1971, p.4938512.
- R. Van Slyke and H. Frenk, "Network Reliability Analysis; Pert I," Networks, Val. 1, No. 3, 1971, p.279-298.
- R. Van Siyke and M. Frank, "Reliability of ComputersCommunication Networks," Proceedings of the 1971 Winter ACM Simulation Conference, December 1971, p.71=62.
- H. Frank and I. T. Frisch, "The Design of Largesteate Networks," Proceedings of the IEEE, Vol. 60, No. 1, January 1972, p.6-11.
- Ha Frank, R. E. Kahn, and L. Klainrock, "Computer Communication Network Design Experience with Theory and Practice." Proceedings SJCC 1972, p.255-270; also in Networks, Vol. 2, No. 2,1972, p.155-166.
- W. Chou and M. Frank, "Routing Strategies for Computer Network Design." Proceedings of the Migroweve Research Institute International Symposium on Computer-Communications Natworks and Teletreffic, April 1972, p.381-389.
- A. Kershenbaum and R. Van Slyke, "Computing Minimum Spanning Trees Efficiently," Proceedings of the ACM Annual Conference, August 1972, p. 516-527.
- He Frenk and We Chous "Topological Optimization of Computer Networks," Proceedings of the IEEE, Vol. 68, No. 11, November 1972, p. 1385-1397.
- H, Frank, "Optimal Design of Computer Networks," in Computer Networks, R. Rustin (Ed.), Prentise-Hall, 1972, p.167-185.
- Mg Gerle, "Deterministic and Adeptive Routing Policies in Pecket-Switched Computer Networks," ACM Third Date Communications Symposium, Florida, November 1975, pg23m26.
- W. Chou, H. Frank, R. Van Slyke, "Simulation of Cantrolised Computer Communications Systems," ACM Third Data Communications Symposium, Florida, November 1973, p. 121=138.

- W. Chou and A. Kershenbaum, "A Unified Algerithm for Designing Multidrom Teleprocessing Networks," ACM Third Data Communications Symposium, Florida, November 1975, p.148=156.
- H. Frank, "Providing Reliable Networks with Unhaliable Components." ACM Third Date Communications Symposium, Floride, November 1973, p. 161-164.
- Is Gitman, R. M. Van Slyke, and H. Frank, "On Splitting Rendom Accessed Broadcast Communication Channels," Presendings of the Saventh Hawaii International Conference on System Sciences, January 1974, Sp.

1.4.1.3 The Telephone Compenies

Building a nationwide communications network means doing business with a number of telephone companies; further, if the network is to have oversees or foreign components, one must also deal with verious international services and national feet Telephone Telephone. As mentioned in a previous section, there exists within the U.S. government an agency called DECCG, which other parts of the government can ask to produce communications services. ARPA used DECCG to produce most, if not all, of the demostic circuits for the ARPANET. In some special cases, for example, the oversees lines to Hawaii and Europe, ARPA dealt directly with the relevant communications euthorities.

The precedure for using DECCO to produce communications service is to send them a Telecommunications Service Request or TSR. For instance, to obtain a circuit from UCLA to RAND, ARPA would send a TSR requesting the necessary circuit, and DECCO would then negetiate with the relevant telephone companies and obtain the appetited service at the best price. In the case of a circuit from UCLA to RAND, for example, most likely the service would be produced from General Telephone, the deminant telephone company in the Los Angeles area.

In the exemple just qiven, a requested of redit fell completely within the jurisdiction of a single telephone company, to handle instances which the requested service spanned the lurisdictions of more than one telephone company, the Bell System set up another sempeny, selled ATST Long Lines, and in the case of many of the ARPANET sirguits, the company from which DECCO procured the service was Long Lines, which in turn presured the components recessary to make up the requested service from the regional telephone companies.

For the first several years of the ARPANET development, the Lend Lines sustemer representative to ARPA was Mr. Ken Stanley. A dustomer representative's job is to make the sustemer aware of the kinds of service available and to keep him happy with the service he receives. Fertunately fer the ARPANET, the Bail System understood that a customer building a nationwide network needed the assistance of some central individual with a broad drasp of the requirements rather than having to rely on the usual miscellaneous dealings with lossit telephone compenies. Thus Stanley, who from his position in Long Lines was already in contact with his counterparts in the many regional telephone companies, was parmitted to use his network of contacts to provide informal coordination of the antire Bail System service to the ARPANET.

For instance, installing a telephone girquit between two IMPs requires that the IMPs and the telephone companies at both ends all be ready simultaneously. It is useless for the IMP supplier and one of the telephone companies to strain to make a saheduled date if the other telephone company cannot make that date. Similarly, it is useless for the telephone companies to strain to make a date if the IMP supplier cannot. Stanley took it upon himself to coordinate all such interdependent events. By wirtue of the Bell System's willingness to provide this critical coordination, an extremely smooth and efficient relationship has been built up between the IMP suppliers, ARPA, and the member companies of the Bell Systems. For a network of the size and domplexity of the ARPANET, there is supplieringly little trouble with the precurement and operation of the telephone circuits.

1.4.1.4 Bolt Berenek and Newman Inc.

The centrest to construct the IMP for the ARPANET was awarded to Selt Beranek and Newman Inc. (BBN) of Cambridge, Massachusetts, where it was cerried out in a group under the leadership of Mr. Frank Meart. Once the ifrat IMPs were installed, BBN continued to play a central role in the evolution of the network, operating it and maintaining it as well as doing the development necessary for several major enhancements of its capability.

In addition to Frank Heart, a number of other individuals at BBN have been invalved with the development of the ARPANET. The names of many of these individuals may be found as euthors of the papers on the ARPANET which have some out of BBNs. However, one individual, Rebert Kehn, stands out. After severely ears of participation in the group working on the network at BBNs, he moved to the IPT office at ARPA from which he has probably denoted to promote and support the tentinued advance of packet-switching technology than any other individual with the possible exception of Reberts himself.

A bibliography of ARPANEToreleted reports and papers written by members of the BBN staff to included as Appendix A.

1.4.1.5 The Network Intermetten Center

The messelbility of distributed resources service with it the need for an information service (either sentralized or distributed) that enables users to learn about those resources: This was resourced at the PI meeting in Mishigan in the sering of 1967. At the time, Doug Engelbert and his group at the Stanford Research Institute were sireedy involved in research and development to provide a semputer-based facility to sugment human interaction. Thus, it was decided that Stanford Research Institute would be a suitable place for a "Network Information Center" (NIC) to be established for the ARPANET, With the beginning of implementation of the network in 1969, sonstruction also began on the NIC at SRI.

The NIC provided several services. It maintained a list of network merticipents and distribution lists for various special interest groups within the network semmunity. An archive of various document series was maintained. Decuments sould be sent to the NIC with instructions for duplication and distribution to the membership or one or more of the special interest groups. A highly structured data base senetruction, manipulation, and display system (salied NLS) was made eveilable oneline for use

over the network. A list was kept of the resources eveilable on hosts throughout the network. The various ARPANET protocol specifications were maintained oneline at the NIC.

The NIC has had a hand in the production or distribution of hundreds and hundreds of documents related to the ARPANET. A few of these decuments describe the activities of the NIC (tself or apa otherwise of special interest within the ARPANET) a bibliography of these follows.

8fb1fography

- O. C. Engelbert, *Network Information Center and Computer Augmented Teem Interestion*, Stanford Research Institute, Augmentation Research Center, January 30, 1971, 185p.
- D. C. Endelbart, "Oneline Team Environment", Stanford Research Institute, Augmentation Research Center, June 2, 1972, 272p.
- D. C. Endelbart, "Coordinated Information Service for a Discipling or Mission-oriented Community", Tipogapariza, Rest. Resease, and Euture. ... the said Landon Community Community. Community Commu

Current Catalog of the NIC Collection, NIC No. 5145.

Current Directory of Network Participants, NIC No. 5158.

Network Resource Handbook, NIC No. 6746.

Current Network Protocols, NIC No. 7184.

1.4.1.4 The Network Measurement Center

Kleinrock sontinued his work as a member of the faculty of UCLA and it was quite natural for him to participate in the study and specification phase of the ARPANET development in the two years that preceded 1970. By the time the IMP RFG was released, it was given to ARPA that Kleinrock and his group at UCLA were the garrest people to set up and staff a Network Measurement Center (NMC) for the ARPANET. Plans were therefore laid to have UCLA be the site of the first IMP installation, to allow early connection of an SDS SIGMA 7 host which was to be used to support the NMC tasks. The NMC had the responsibility for much of the

analysis and simulation of the ARPANET performance, as well as direct measurements based on statistics gathered by the IMP program. While Kleinrock himself was the guiding force at the NMC, over the years he had a series of students or staff members who supervised the day-to-day measurement work, including Gerald Cole, Vinter Cerf, Holger Opderbeck, and William Neylor (each obtained a UCLA Ph.D., although not necessarity for their NMC work).

There has been a series of doctoral dispertations written by students at the UCLA School of Engineering and related to the work of the NMG. A list of of these theses, as well as other relevant publications by the faculty and students associated with the NMC, is included in the following bibliography.

Bibliographyt

L. Kieinreck, "Communication Nets: Stechastic Message Flew and Delay," Dover Publications, New York, 1964, 289s.

is Klainpack, "Models for Computer Natworks," Proceedings of the International Communications Conference, June 1969, p.21, 9021, 16.

L. Kleinrock, "Comparison of Solution Methods for Computer Network Madels," Presentings of the Computers and Communications Conference, October 1969.

^{*}I have asked Len Kieinreck to suggest additions and deletions from this bibliography we DCW,

- L. Kieinrock, "Analytic and Simulation Methods in Computer Natwork Design," Proceedings SJCC 1970, p.369-579.
- L. Kieinrock and G. L. Fuitz, "Adaptive Routing Techniques for Store-end-Forward Computer Communication Networks," International State of the Art Report No. 6; Computer Networks, Infetein Information Ltd. Meidenhead, Serkshire, England, 1971, p.341-362; also in Proceedings of the International Conference on Communications, Mentagel, Guebec, Canada, June 1971, p.30.1-39.8.
- G. Cole, **Perfermence Measurements on the ARPA Computer Network; **
 Proceedings of the Second ACM/IEEE Symposium on Problems in the Optimization of Data Communications Systems, October 1971; p.39-45; also in IEEE Trans. on Communications, Vel COM-28, No. 3, June 1972, ph38-636.
- G. D. Cole and L. Kleinrock, "An Analysis of the Separation between Packets in a Storemendereward Network," Proceedings of the Microwave Research Institute International Symposium on Computer-Communications Networks and Teletraffic, April 1972, p. 1-3.
- V. G. Cerf and W. E. Neylor, "Storede Considerations in StoremendeForward Message Switching," Paper presented at HESCON Conference, September, 1972, Sp.
- V. G. Cart and H. E. Naylor, "Selegted ARPA Network Messurement Experiments," Proceedings COMPCON 1972, September 1972,
- L. Kleinrock, "Computer Networks," in Computer Science, A. F. Cordenes et, el. (Edg.), John Wijey and Sons, Inc. New York, 1972, p.241-254.
- L. Kielnrack, "Performance Models and Measurement of the ARPA Computer Network," Precedings of the Internetional Symposium on the Design and Application of Interestive Computer Systems, Brunet University, Uxbridge, England, May 1972, 38p.
- L. Kleinrock, "Burvey of Analytical Methods in Queueing Networks," in Computer Networks, R. Rustin (Ed.), Prentice=Hell, 1972, p.163=205.
- D: Cantor and M: Geria: "The Optime! Routing of Messages in e Computer Network Via Mathematical Programming:" Programdings COMPCON 72, September 1972: p:167-178;

- L. Fratta, M. Geria, and L. Kieinrock, "The Flow Deviation Methods An Approach to Store-and-Forward Communication Network Design," Networks, Vol. 3, No. 2, 1973, p. 97-134.
- L. Kleinrock and S. S. Lem, "Pecketmewitching in a Slotted Satellite Chennel," Proceedings NCC+E 1973, p.783-710.
- L. Kleinrock, "Queyeing Systems; Theory and Application," Wiley Interacience, New York, 1973.
- L. Kleinrock and S. Lem, "On Stability of Packet Switching in a Randem Multi-Access Broadcast Channel," Presendings of the Seventh Haweii International Conference on System Sciences, January 1974, 40.
- D. G. Centor and M. Gerle, "Capacity Allegation in Distributed Computer Networks," Proceedings of the Seventh Hawaii International Conference on System Sciences, January 1974, p. 115-117.
- G. Cole, "Computer Networks Measurements: Techniques and Experiments," PhD Thesis, UCLA-ENG-716E, Computer Science Department, School of Engineering and Applied Science, UCLA, October 1971, 358p.
- G. L. Fults, "Adaptive Routing Teghniques for Message Switching Computers Communications Networks," PhD Thesis, UCLA-ENG-7252, Computer Science Department, School of Engineering and Applied Science, UCLA, July 1972, 418p.
- M_e Gerla, "The Design of Storemendeforward (\$/f) Networks for Computer Communications," PhD Theeis, UCLA-ENG-7319, Computer Seience Department, School of Engineering and Applied Science, UCLA, January 1973, 399p.
- L. Kleinrock, "Computer Network Research," UCLA February 1978, 26p.
- L. Kleinreck, *Computer Network Research, * UCLA August 1979; 189p.
- L. Kleinrock, "Computer Network Research," UCLA June 1971, 138p.

L. Kleinrock, "Computer Network Research," UCLA December 1971, 115p.

L. Kleinreck, "Computer Network Research," UCLA June 1972, 113pg L. Kleinrock, "Computer Network Research," UCLA Desember 1972, 78pg

L. Kieinrock, "Computer Network Research," UCLA June 1973, 64p.

1,4,1,7 The Network Werking Group and Heat Involvement

The initial design of the ARPANET as senteined in the REG went some way toward specifying the protocols for inter-IMP communications and for IMP-to-host communication. Lass explicit attention was given to host-to-host communication, this area being left for host sites to work out among themselves.

To provide the hosts with a little impatus to week on the host-to-host problems, ARPA essigned Elmer Shapiro of SRI "to make something happen", a typically vague ARPA assignment: Shapiro esliad a meeting in the summer of 1966 which was attended by programmers from several of the first to hosts be connected to the network. Individuals who were present have said that it was glear from the meeting that at that time, no one had even any very glear notions of what the fundamental hosts-to-host protocol issues might be.

After the Sante Berbera meeting, S. Crocker, S. Carr, and J. Rullfann mee again in the summer and fall of 1968 to continue discussion of hostetechost protocol issues. Their early thinking was at a very high level, e.g., the feesibility of greating a portable front-end package which could be written once and moved to all network hosts; a host desiring to send data to another

host would first send a date description to the receiving host how to interpret date coming from the sending host. On Vetentine's Day, 1969, the first meeting of host representatives and representatives from the NMC and NAC, slong with the IMP contractor, was held at BBN in the middle of an enormous snew storm.

The several individuals considering heateteahest protocol issues had never estually received an explicit charter from ARPA to do this work, and they worried that someone might think they were acting out of place. For this reason, in April 1969, they established a series of working notes estled Request for Comments (RFCs), which could be effectived to let others knew what they were doing and to obtain the reactions and involvement of other interested parties. They called themselves the Network Working Group (NWG).

The NWG eventually grew quite large, with representatives from almost every host site in the network participating, and mountains of paper was circulated describing and commenting on various protocols. There were also occasional mass meetings. From about the time it was decided that he would go to ARPA until

near the time he left ARPA, Stephen Crocker gerved as chairmen of the NHG. By the beginning of 1972 the NHG had grown too large, but much of its work was done as large numbers of hosts were sommunicating over the network. From this point enwand, meetings were limited to those of an executive protocol gammittee which met to discuss general protocol issues and provide guidance for Crocker, and to those of various subsommittees, e.g., the group interested in the Remote Job Entry pretocol. Even after the big meetings ecopped, mast participant working notes were directed to most other participants in the network.

Greduelly the setivities of the NWG began to diminish. Many of the heat site personnel who had originally been satius moved on to other tasks, and new users joining the network tended to use the defined protogels rather than becoming involved in their specification. As Crocker's time for the NWG group become ingressingly limited, he appointed Alex McKensie and Jon Postel to serve jointly in his piece. McKensie and Postel interpreted their task to be one of sedification and coordination primarily, and after a few more spurts of estivity the protosel definition process settled for the most part into the status of a maintenance effort. Teday when one of the engeneentral figures in the host protogol design process girculates an RFC among his

DRAFT

remaining contemporaries, it may be more for all times" sake than for any other resson.

Further activities of the Network working Group will be described in the section below on the development at host-to-host protocol.

1.4.2 Initial Subnet Decion

SSNFR proposel for the IMP was for the most part compliant with the requirements of the RFG. In some important respects, however, the BSN IMP design diverged from those requirements or added constraints that were not in the RFG.

The BBN design took the idea of inserting a communications processor between the host and the network to a logical extremet it specified that the IMP be used only for communications functions, and that there be meximum legical separation between IMP and host. Thus, the design did not provide for IMP repeaking of binary messages on behalf of a host on for doing character conversion for a host. Further, the design precluded all host programming of the IMP and any sentrol of the IMP by regular hosts. This pertion of the IMP design undeubtedly was correct and contributed to the repid development of the subnetwork of IMPs.

The initial subnet design also specified minimal control messages between an IMP and its hosts, many fewer than envisioned in the RFP. As the network later developed, it proved useful to additional such control messages.

The IMP design estind for a hardened IMP, which would be robust in the face of physical and electrical abuse. While statistics have shown that hardening did help, it was eventually desided not be be worth the added price of the hardware.

The depablity to loop all interfaces was included in the design. This has proved to be of enormous operational importance.

A decision was made to reload IMPs initially from paper tape, and each IMP was provided with a paper tape reader. It was always intended that this progedure would sooner or later be replaced by loading through the network, and it was. The use of the paper tape reader at each IMP was probably a useful simplication in the beginning.

Program debugging was also specified to take place from a local terminal, in the interests of simplicity, with the knowledge that cross-network debugging could be added later. A cross-network debugging sepablity was added even before the first IMP was delivered.

A delivery schedule of one IMP per month in months eight through eleven after contract start, instead of all four at month nine, was assumed.

In keeping with the strict independence of host and IMP, the IMP was to gether statisties routinely and transmit them periodically to specified hosts rether than permitting hosts to control the IMPer statistica-taking gapabilities. This proved satisfactory.

The initial IMP design was responsive to the RFQ in one particular which was changed at the first meeting between BBN, ARPA, and hest representatives. RFQ had called for one hest per IMP and aix directs per IMP, although it also asked that more hosts per IMP be an option. Almost immediately it became elegates many hest sites had more than one host to connect to a single IMP. Thus before the first IMP was delivered, the design was changed to permit two, three, or four hosts on an IMP elong with five, four, or three lines. This was a simple change to implement.

In the area of the host-to-IMP protectly, the initial IMP design specified this protectl as required. Unfortunately, seme aspects of the host-to-IMP protectly had eignificant detrimental effect on the design and performance of the other protocols. Particularly unfortunate have been the adknowledgment system, the retransmission system, and the message identification system initially suggested by the host-to-IMP protocol.

It is crucial for the IMPs to limit the rate of flow of host traffic into the net to the rate at which that traffic is being taken but, in order to prevent subnetwork condestion. The IMPs first attempt, which was insufficient, took the form of suggesting that each message in a conversation should be held by the sending host until an endatowend asknowledgment for the previous message was received. This suggestion was adopted as part of hestotochost protocol upon which all the higher standard protocols are based. As a consequence, the bandwidth of a single host-to-host protessi connection is severely limited; given the ARPANET response time using 50 Kbs lines, waiting for a destination-to-source ecknowledgment between messages typically limits connection bendwidth to about 18 Kbs, in contract to the 48 Kbs possible with a constant streem of messages,

It was originally thought that the ARPANET would lose a message so selden that there was no point in hosts ever bothering with message retransmission. Unfortunately, reselving various possible legiums has required the subnetwork to distard a message occasionally, and the topology of the network has evelved into long series of machines and lines that increase the probability of involuntary message less. Hewever, the hostetowhost protogol followed the initial thought and did not provide for message

retrensmission. Given the realities of the probability of message loss in the network and given the hesteto-host protocol, which is inerdinately sensitive to any abnormality, the host-to-host protocol (and protocols besed on (t) has proved quite unreliable, and the original host-to-IMP protocol must be held partly responsible.

part of the hostetowIMP protocal endetowend asknowledgment system described above, the hostetowIMP protocol apecified an 8-bit message identification number and suggested that all messages in a single conversation corry this same identification numbers in fact, messages with identification numbers were not guaranteed to be delivered in the order sent. Eight bits is probably insufficient to identify uniquely (for the purpose of possibly required retransmissions) outstanding messages when successive messages in a conversation are sent without weiting for an endetowand agknowledgment. Thus the small identifier prevented reliable Aubit Message high-bandwidth connections.

The IMP/heat proteco; has been shanged so that it is no longer necessary to wait for the endetowend acknowledgment; massage order is now preserved except for priority

DRAFT

Section 1

considerations; cases requiring message retransmission are unambiguously reported to the sending heat; and the message identifier has been expended to a sufficient size; Unfortunately, there is great inartia in the heatwoodst protocol due to the large number of implementations that would have to be changed. The host-to-host protocol is still incapable of simultaneous high performance and perfect reliability.

1,4,3 Subset Development

The first four IMPs were developed and installed on schedule by the end of 1964. No sooner were these IMPs in the field then it become clear that some provision was needed to connect hosts relatively distant from an IMP (i.e., up to 2006 feet instead of the expected 56 feet). Thus in early 1976 a "distant" IMP/host interface was developed. Augmented simply by heftier line drivers, these distant interfaces made alsor for the first time the falledy in the assumption that had been made that no error control was needed on the host/IMP interface because there would be no errors on such a local dennection.

By mid-year of 1970, a series of network performance tests were being serried out. These unsovered some flaws which were quickly corrected, and some problems which looked more worrisome. Also by mid-year, a rudimentary version of the network centrol center was established at BSN.

As the year were on, eitem continued to be added to the network, the IMP program continued to be improved, NCC development continued, the first 238.4 Kb sircuit was tested between two IMPs, and design for a verson of the IMP able to support direct terminal connection was begun. The latter was dailed the terminal IMP.

By early in 1971 melor problems with the IMP flow control and storage allocation techniques were found, and design began an improved algorithms. It is interesting that ence these problems were found, and they were serious enough to sempletely halt network operation, the network continued to give adequate service for many, many months while improvements were designed and implemented. The hosts were simply asked to not use the network in the way that saused the subnetwork problems, and the heats did as they were esked.

About three-quarters of the way through 1971 the first two TIPs were delivered, providing ARPANET access for the first time to users without their own heats or access to terminals on same other organization.

By the beginning of 1972 it was recognized that even the distant version of the IMP/host interface was not sufficient, and design for a IMP/host interface for use over communications einquits was begun. The evolution of the IMP/host interface is worth a little additional comment. The initial bitwser(a), saynghronous, nonvercementabled IMP/host interface was essentially specified in the RFQ, in an effort to simplify network connection for the hosts. This nonestandard interface

may have been of some benefit in simplifying the host gonnection; However, its greatest virtue was the separation it put between the IMP and the hast. The IMP and host did not have to worry about each other's word size, and they did not have to worry about each other's timing constraints. It seems likely that having to worry about these issues would have deleyed network operation. However, this interfede also resulted in a hodge-podge of interfage variations, each designed for more distant operation then its predecessors, and none except the first was vary elegant. For any new network, which need not fear the now preven packets-witching technology, it would sleerly be better to use an industry standard communications interfage, e.g., HDLC, for every IMP/hest connection.

In the first helf of 1972 the TIPFe sepability was expended to support TIPstosTIP magnetic tape transfers. While this option was successfully used between two network sites, it was never very elegant, and simest immediately had to be partially redons. Also, a messive change in the IMP software was undertaken to correct the previously discovered flow control and storage allocation problems. In the second half of the year, the new version of the IMP program was released in many small increments, and the design of a new, ten times more powerful IMP was begung

The beginning of 1973 brought the first setellite link in the network, from California to Hawaii. Also, with network traffic repidly increasing, a number of subnet reliability problems developed which had to be corrected. By mideyear, a pelr of TIPs had been shipped to Europe, for use in Norway and These brought pienty of operational troubles. For the first time, elecuits had to be obtained from a foreign PTT, the pircuits were relatively slow at 9.6 Kb, and like Hawait, these TIPs were on a long spur off the network rather than being connected doubly as IMPs typically are. As 1973 continued, nodes continued to be delivered, but there began to be a lew level of switching of mode logations, to optimize the use of various. IMP configurations and as sites same on and went off the network, Certain improvements were also made to correct problems with the routing elgorithm. As 1973 ended the first very distant hosts were connected to the network.

In 1974 there were major efforts to make the network more operational. Subnatuork reliability was improved as wes TIP-to-host communication reliability. Methods for providing TIP access dontrol and accounting and partitioning of logical subnatuorks of hosts were developed. Methods were developed to selectively reload sections of IMP memory.

In 1975 network development glowed up and the network took on more and more of an operational eppearance. Major network developments in 1975 (noluded delivery of the first Pluribus IMP, modification of the IMP and TIP software to support more than sixty-three IMPs in the network and attachment of the first two Satallite IMPs to the edge of the network, By the end of 1975 the network was under OCA management.

Looking back, the subnet development between 1969 and 1975 appears relatively smooth, although there were many times during that period when those intimetely involved felt they were trying to solve one sridis or another. The network grew slowly anough, and the basic technology and implementation was flexible and robust enough, that many problems, both major and minor, which naturally propped up with this new development were for the most part corrected before they obstructed the work of too many users. The fact that the network was explicitly an experiment no doubt also made users mere tolerant.

1.4.4 Host Protogol Development

Specifications exist for the physical and logical message transfer between a host and its IMP. These specifications are denerally delied the IMPetowhest protocol. This protocol is not sufficient by (teelf, however, to specify the methods communication between processes running in two possibly dissimilar hosts. Rather, the processes must have some agreement as to the method of initiating sommunication, the interpretation of transmitted data, and so forth. Although it would be possible for such agreements to be reached by each pair of hosts (or processes) interested in communication, a more general arrangement is desirable in order to minimize the amount of implementation nesessary for networkswide communication. Accordingly, the heat organizations formed a group (the Network Working Group or NWG, introduced above) to feetilitate an exchange of ideas and to formulate additional specifications host-to-host sommunications,

The NWG adopted a "layered" approach to the specification of dommunications protocols, wherein the higher layers of protocol use the services of lower layers; the adventages and disadventages of the layered approach are discussed elewhere in

DRAFT History

Section 1

this report. As shown in Figure 3, the lowest layer is the

Figure 3 == Layered Reletionship of the ARPANET Pretocols

IMPromhest protocol. The next layer (called the host-to-host layer in the figure) specifies methods of astablishing communications paths between hosts, menaging buffer space at each end of a communications path, ets. Next, the Initial Connection Protocol or ICP specifies a standard way for a remote user (or process) to attract the attention of a natural host, preparatory to using that heat, The ICP provides the analog of the user pressing the attention button at a local terminal on a host. In the next layer is the Telecommunications Network or TELNET

protocol, which was designed to support terminal access to remote hosts. TELNET is a specification for a network standard terminal and the protocol for communicating between this standard terminal and a host. The next logical protocol layer consists of function oriented protocols, two of which, file Transfer Protocol (FTP) and Remote Job Entry protocol (RJE), are shown in the figure. Finally, at any point in the layering process, it is possible to superimpose adiasa protocols.

In the following subsections we discuss in some detail the events in the evolution of the host-technol and TELNET protocols, and the events in the evolution of a number of other protocols in somewhat less detail. Appendix 8 contains a chronology, supplied by Jon Possel, of the host protocol development experience.

_ Î

1.4.4.1 Host-to-Host Protocol

The Network Working Group was established in early 1969. By December 1969 an initial host-to-host protocol had been specified which supported communication between a terminal on one host and a process on another host. At a meeting in Salt Lake City in December 1969, the initial protocol specification was described to Lawrence Roberts of ARPA who was unhappy with it because the initial plan would not support transmission of electronic mediover the network. He instructed the Network Working Group to "go back and get it right."

By the apring of 1978 several suggestive versions of a host-to-host protocol had been developed, and a relatively formal meeting of the NWG was held at UCLA before mid-year at which the letest version of the protocol was described. Resetions to the described protocol were very negative. In June of 1978 there was a series of meetings held at UCLA and Harverd at which people from these two institutions tried finally to settle upon a host-to-host protocol and specify how it should be implemented. In August of 1978 some of the more general (and some thought more exotic) aspects of the host-to-host protocol being sensidered were ordered dropped from the protocol by Barry Wessler of ARPA,

thus administratively elearing away some of those (saues which had prevented agreement. The NWG discussion continued at the 1970 Spring Joint Computer Conference; in particular, there was discussion between Crocker and Roberts regarding the formality to be sought for the protecol, and ARPA approvals required, and so forth. Another NWG meeting was held at the Fell Joint Computer Conference in November 1976 in Heuston, Texas.

By seriy 1971 ARPA was growing increasingly impetient with the progress being made to settle on a heatesembest protect. At a NHG meeting held in mideFebruary 1971 at the University of Illinois, a subcommittee was appointed to leak at the host-to-host pretectl to see what changes were immediately desirable or negaciary. This subcommittee want directly from Illinois to Cambridge, Massechusetts, where it met for two days, wrote an interim report, and then reconvened a menth later in Las Angeles. It appears that with the efforts of this semmittee (known as the Mestatewhost protocol glitch elsening committee) the design of the ARPANET host-to-host protocol was finally coming close to being settleds

At about this same time ARPA was beginning to exert great pressure not only to get the hest-to-host protocol settled but

also to get it implemented by the hosts. Some heats had already implemented an early version of the protocols in feat, at least one host, with especially productive programmers, was penalized for its efficiency by implementing several suggestive versions of hoststochest protocol. At a NWG meeting at the Spring Joint Computer Conference in Atlantic City in May 1971, Roberts strongly admonished the NWG to "finish something." At that NWG meeting Alex McKenzie took on the teak of writing a definitive specification of the hoststochest protocol so not to invent new protocol, but to write down what had been degiced.

In October 1971 the final big NWG meeting was held at M.I.T., and was preceded by a programmers workshop at which differences in implementations were clarified and climinated. In January 1972 a McKenzie document describing the protocol was published and the ARPANET host-to-host pretocol has remained essentially unchanged singe.

1,4,4,2 The Evolution of Teinet

Early in the development of the ARPANET it became glear that a major function of the network would be to provide remote use of interactive systems. To allow a user at a terminal (connected to his lodel host) to control and use a process in a remote host, as iff he were a local user of that remete host, a special mechanism was required. The problems to be eversome are tegions example, the typical host expects its interactive terminals to be physically ettached to the individual parts of its hardware terminal scanner rather then indicativ attached via a multiplexed connection to the network; a given host expects to communicate with terminals with sertain characteristics (e.g., oniv half-duplex, line-etwestime, physical acho, ESCDIC charecter set, 134.5 baud) while a remote user's terminal might have gempletely different . abarestaristics. (0.0.. - fulleduplex. sharacter-at-cotime, no character acho, ABCII character set, boud). The TELNET protocol was an attempt to provide the special mechanism necessary to permit such communication.

As early as 1969 a few hosts had been programmed on an ad the basis to permit terminal access from another host. In 1971 an NWG subcommittee was formed to consider the general problem of

supporting intersetive use of applicant hosts by users at applicant remote terminals. There was great dentroversy in the sommittee discussions, focusing on four issues: character set, connection establishment, echoing, and interrupt capability. By late 1972 there was anough consensus so that widespread implementation of an early version of the TELNET protocol had been accomplished.

Despite widespreed implementation of the early TELNET protocol, its heavy and effective use, and numerous attempts to desire it complete, discussion of it continued. There were several problems with the early version:

- in Despite the attempt to permit a minimal implementation well suited to the constraints of small hosts, there was no well-defined minimal implementation. Even if some TELNET feature was not desired for a given implementation, it had to be provided in same come other implementation commended its use.
- 2. The control atructure was inedequate. For example, unless some exceedingly constraining assumptions were made, it was possible for the two ends of a TELNET connection to loop, commanding each other to take opposite actions.

- In asymmetry of TELNET connections precluded one end from initiating certain functions, such as echolog behavior. This aeriously constrained the use of TELNET protocol for character communication between processes not serving terminals, a role for which it would otherwise have been well suited and for which it was elready frequently used in the absence of any better protocol.
- 4. The impus of interfacing sharesterwateautime hosts to Timewateautime hosts was poorly handled.

8y early 1973 it had become apparent that minor adjustments to the early IELNET protocol would not solve these problems and that some fundamental abendes were needed. A new subcommittee mat and, with the previous experience to guide them, developed several fundamental principles. These new principles, when added to the service principles of the Network Virtual Terminal and the remote interrupt (synch) mechanism, resulted in a revised TELNET protocol which selved meet of the earlier problems that had preciuded universal seceptance of the pretocol.

There was such enthusiesm for the new vereion that a schedule for "repid" (within the year) implementation was laid

However, the implementation of the new TELNET protocol out. proceeded mere slowly than expested. Only in the past year implementations been widely available. In retrospest, there were several reasons for the delay in the implementation: 1) at the time the revised implementation was protocal scheduled. implementation of the initial version had been completed and host system managers had not budgeted resources for a segond implementations 2) about this time ARPAPs research interest in the network was dealining and the network was entering a period of status que operations 3) despite initial belief that a cleen method of phesing over from the initial protocol to the revised protocol existed, none was found by most implementars and sonsequently mest chose to provide a complete implementation of the revised protecol to exercte in perellet with the initial protocol: and 4) implementation for the most prevalent user host, the TIP, proved to be very diffiguit (because of the TIPFS limited memory) and time-gontuming, thus implicitly relieving pressure on the server hosts to implement the revised protecti.

At the time this is being written, in late 1977, the new TELNET protecol has been the accepted standard for several years, and it is widely implemented and used.

1,4,4,3 The Evalution of the Other Host Protocols

There are several other host protects the evalution of which should be briefly mentioned.

The File Transfer Protocol started out as two protocols, a Deta Transfer Protocol and a File Transfer Protocol. To over-simplify, the Data Transfer Protocol was to specify the formet of data being transferred and the File Transfer Protocol was to specify how it was transferred, Eventually, the File Transfer Protocol alone was defined with a data partial and a control portion. After the final push to specify the FTP, relatively little additional work was done, cansisting only of a little affort to clean up fundamental aspects of the protocol, and a good bit of work reconciling the Preply codes that different hosts used to indicate FTP-related events.

Before a Remote Job Entry protocol could be defined by the NMG as a whole, UCLA's ISM 368/91 host, a batch oriented host, needed seme AJE=1ike protocol with which to serve a few users who wented early assess to the semputing power of that particular host. Thus, led by the UCLA group, a protocol called the Remote Job Service or AJS protocol was defined and implemented. The NMG eventually got around to working on the problem of a Remote Job

Entry protocol and undertook a relatively massive effort to define such a protocol. However, by the time the RJE protocol definition was finished, half a dozen or so hopes had already implemented to interim RJS protocol. Since these included most of the hosts on the network interested in supporting remote batch, there was little insentive for them to implement the new RJE protocol. Thus, today the RJE protocol is serefully specified but to our knewledge is not implemented enywhere, and the RJS protocol proveits.

Another protocol much discussed within the NWG for a time is the Date Reconfiguration Service protocol. The Date Reconfiguration Service was to be a facility which would come to have tables to convert from the format of any host on the network to the format of any other. For ressons which are unglear, the service was never implemented.

Moving upward in sophistication, another protocol that was the subject of early discussion was one for graphics. Several versions of a graphics protocol were specified but there was never widespread implementation of any of them.

In addition to the heat-to-heat protocol which was finally specified after much iteration, a number of alternative protocols

were suggested by verious members of the NWG. Before the host-to-host protocol was settled upon, Richard Kaiine and David Welden each suggested an alternative protocol. Even after the adoption of the host-to-host protocol, there was some discussion of experiments with a protocol derived from the Welden suggestion. More recently, as part of the ARPA-ponsored National Seftware works project, 177, Richard Schentz, and Robert Thomas have designed and implemented a host-to-host protocol known as MSG. Another protocol, known as TCP, deserves special mention.

Near the time of the formation of the International Network Morking Group, as network interconnection began to be of great interest, discussions began on a standard interenetwork protocol, particularly one which would correct seme of the shortdomings of the ARPANET host-te-host protocol, At the AFIPS 1973 NGC in New York City a meeting was held at which sertein ideas for a new host-te-host protocol were discussed. After some additional correspondence, Robert Kehn of ARPA and Vinton Corf, then of Stanford, get tesether and designed a protocol known as TCP. Other members of INMG, perhaps not satisfied that TCP represented an international standard, continued developing still snother host-te-host protocol (Cerf also participated in this later

effort). TCP quickly become ARPA's choice of the host-to-host protocol to be used in situations where the ARPANET host-to-host protocol was insufficient or where inter-nativarying was required; with ARPA support, several TCP implementations were done and the protocol has some into relatively widespread use within the ARPANET, and its use is still spreading. It appears that ARPA envisions the day when TCP will replace the ARPANET host-to-host protocol. Meanwhile the host-to-host protocol, Meanwhile the host-to-host protocol that the rest of INWG was working on was finished, and documented, just as the PTTs and North American common carriers submitted the X,25 standard to CCITT; so the INWG consensus protocol will most likely play little operations; role in the ARPANET or elsewhere,

DRAFT

1.4.5 Network Growth we A Summary

In the following three subsections we consider three aspects of the growth of the network; the traffic growth, the growth of the network topology, and the increase in the number and type of hosts on the network.

1.4.5.1 Traffic Growth

In early 1973, Roberts presented a surve of everage host internade traffic growth for the networks which showed the level of internade network traffic to be increasing at a rate of a factor of ten every ten months. Internade traffic means traffic sent from a host on one nade to a host on a different nade; is earlit does not include traffic sent between hosts on the same nade; because on this rapid nate of growth, Roberts predicted the network would run out of papasity in nine months. As shown in the following figure, shortly after Roberts prediction the nate of internade traffic growth degreesed sharply to roughly a factor of two every twenty menths. It is interesting to speculate on the reason for this sharp degreese.

e L.G. Roberts, "Network Rationals: A Sayear Resvaluation," Proceedings COMPCON 1973, February 1973, pp. 3-6.

DRAFT History

Section 1

Figure 4 - Growth in Average Host Interneds Traffic

Before the network existed, ARPA appears to have had a tendency to buy a computer for each of its research contractors. Once the existing computers were put on the network as ARPA appears to have shown greater rejudtance to provide each of its contractors with a host, preferring instead that contractors use existing hosts; in fact, some groups sequired all of their computing over the network. Thus it can be hymathesized that the axisting hosts (and the few new heats that were added) were used remotely more and more and network traffic increased more and more, until the hosts (at least the popular time=sharing hosts) began to run out of sapasity; this made it paintiess for new remote users to attempt to get service. and resulted, turn, in a leveling off of network traffig growth. Therefore, instead of the network running out of capacity as predicted by Roberts, it seems that the hosts ran out of separity while the network still has capacity left.

As already stated, the traffic shown in Roberts? curve and in the figure above includes only intermede traffic. There are two resears for excluding intranede traffic. First, intranede traffic puts a burden on only one nede rather than on the network as a whole. Thus when Reberts, for instance, was attempting to dejouiste the effects of host traffic on network sapecity, he

neturally excluded intranede traffie: Segend, the evaliable intrances traific statistics include some amount of test traffic being looped from a host through its node and back to the same host, and there is no convenient way to separate this looped test traffic from actual data traffic between two hosts on the same node. It is believed, however, that there is accusily a atonificant amount of real traiffs between hosts on the same node. For instance, Kleinrock reports* that during a week-long messurement, the level of intranode traffic amounted to a daily everege of twenty persent of the level of intermode traffie, and in some onecheur intervels the intranode traffic level was as much as eighty percent of the intermode traffic level. A mean of available long term statistiss on inter- and intraneds traffic shows that intranode traffic lavels have everaged between twenty and forty percent of interneds traffic levels. Thus the traffic surve given in the figure above should be sealed up by this factor if all traffic is to be included.

Thet intranede traffic is a significant portion of all network traffic is interesting and probably indicative of four

E. Kielnrock and We Nayler, "On Measured Schooler of the ARPA Network," AFIPS Conference Prospedings, Volume 43, May 1974, pp. 767-780.

phenomena. First, the IMP is a kandy interhest interface, and ende one is installed in a computer center to connect come host onto the network, there is very soon pressure to sonnect other somputers in the somputer senter to the IMP so that desired communication between the computers is possible. Second, when two computers are connected to the same IMP to they may both communicate with other computers in the natwork, communication between the two demputers themselves somes free and begins to happen even if it was not initially thought to be desired. Third, the TIP (e host) has been shosen by several sites as the most flexible available terminal multiplexor and TIPetochost traffic at these sites is likely to be intransde. Fourth, there is a (as yet still weak, but definite) tendency for hosts to be concentrated at a gertain site and therefore eften on the same IMP. The reason for this tendency is that, while some synica would have guessed that every somputer senter manager is trying to build his empire as large as possible, in fast the world of computer center managers appears to instude not only managers whose inclinations are as the cynics sugsed but also many who distike running somputar centers but do so because they need the service supplied by the computer century. Once the network became available, some sites have arranged with some other sites that one site a computer was moved to a second site, and the second site menaged it for the first site which used the semputer over the network via a simple terminal concentrator. A further reason for this tendency (for hosts to be diustored) is the economy of scale possible when only one facility and staff is required for the operation of several computers.

1,4,5,2 Topology

The first ARPANET made was installed at the University of California at Los Angeles in late 1969 and the next three mades were installed in California and Utah.

Pigure 5: The ARPANET in December 1469

By June 1978 three East Coast and two more West Coast nodes were edded, as well as two cross-country lines.

Figure 6: June 1978

IMPs continued to be delivered to the field at an everage rate of approximately one per month, so that by late 1978 there were thirteen IMPs installed in the network. The IMPs were all entirely competible, all being based on the Honeywell 516 computer.

Floure 72 December 1978

DRAFT

Section 1

By two thirds of the way through 1971, two additional 516 IMPs had been installed, the prototype TIP was running at 88N, and two TIPs were exerctional within the network, at MITRE and AMES.

Figure 8: September 1971

The TIPs were based on Honeywell 316 instead of 516 demputers and had as a component a 316-based IMP which was dempletely compatible with the 516 IMP but half as expensive. By early 1972 geveral additional IMPs and TIPs had been installed and the central part of the network between the East and West Coast clusters was beginning to fill out.

Figure 91 March 1972

By August 1972 a third gross-country line had been edded and it was elear that in eddition to the IMPs scattered throughout the center of the country, there were setually clusters of IMPs in four geographic eress, Boston, Washington, D.C., San Francisco, and Los Angeles.

Flaure 10: August 1972

Section 1

There follow onderwearly maps for the years 1973 to 1977 with which the reader can follow the continued growth of the ARPANET topology.

Figure 113 September 1973

Figure 121 June 1974

Figure 13: July 1975

Section 1

Figure 141 July 1976

Figure 15: July 1977

Having shown the growth of the ARPANET topology through a series of geographic maps, it is interesting to consider the evalution of the topology on a more quantitative basis. We use the topologies given in the eleven maps shown above as the basis for the quantitative date shown in the Figure 16 chart. Columns 6 and 7 are missing the first six entries because the HAWAII, NORSAR, and LONDON nades did not exist at the times the six earlier maps were made. The first three entries in columns B through 11 are missing because that information was not kept in the early days of the network. There are several interesting facts that should be noted from the shart. First, the number of network nades has steved about the same since 1975. Despite this, the network throughput has continued to increase. Next, note that intrande throughput actually is a significant fraction of total network throughput. Also, note the peak in everage path length reached in 1974-75, and the subsequent decrease in everage path lengths selected lines were added to the network in 1975=76 es a direct response to network delay problems which occured in 1974-75. Finelly, note that in July 1975 and July 1977, the path from HAWAII to LONDON is equaled by other paths in the networks the NCC (in the Boston area) was 15 hops from both UCBB (Santa Barbare, Californal and FNWC (Monterey, California) in 1975, and

1	1 2 1	1 , 3 ,	1 4 1 5	1 6 1	71 8 1	9 (10	11
DEC69	1 4	2.00	1,33 2			4944		*
JUN70	9	2,22	2,31 4	***	•• •••		₽₹⋟₹	•••
DEC70	13	2.46	2.76 6			= e. ⊕ ■		
3EP71	18	2.44	3,32 7		3,27x	2,892	3,121	6,013
MAR72	23	2,35	5.04 11		4,00x	11,633	21,073	32,706
AUG72	29	2,21	4.65		1,79x	662,502	287,953	970,455
SEP73	40	2,20	5,61 13	9,40	11 3,53%	2,893,138	742.746	3,635,876
JUN74	46	2.17	6.14 13	5,98	12 1,19%	1,125,955	1,513,777	4,639,732
JUL75	57	2,28	6.79 15	6,68	15 .67X	5,179,361	1,918,538	7,097,899
JUL76	58	2.45	5,26,11	5,13	10 .56%	6,627,968	1,961,726	8,589,694
JUL77	 58	2.41	5.37 11	 5. 27	111 .94%	7,051,9221	2,766,054	9,757,976

where the columns contain the following informations

```
Column 1; Date of Map
Column 2; Number of Nodes
Column 3; Average Connectivity
Column 4; Average Path Length
Column 5; Maximum Path Length
Column 6; Average Peth Length, Minus HAWAII, NORSAR, LONDON
Column 7; Maximum Path Length, Minus HAWAII, NORSAR, LONDON
Column 6; Persentage of Node Unavailability
Column 9; Internede Throughput
Column 10; Internede Throughput
Column 11; Sum of Internede and Internede Throughput
```

Figure 16: Some Quantitative Data

Section 1

DRAFT

in 1977 SRIZ (near Sen Francisco, California) is eleven hope from NRL (near Washington, D.C.).

There has been a good deal of actual measurement of the behavior of the ARPANET, and the most detailed discussion of this is presented in the paper entitled "On Measured Behavior of the ARPA Network" (i. Kleinrock and W.E. Naylor, APIPS 1975 Conference Proceedings, Vol. 43, pp. 7674780).

1.4.5.3 Hosts

The first four hosts sonnected to the ARPANET were an SDS SIGMA-7 at UCLA, an SDS-940 at SRI, an IBM 360/73 at UCSB, and a DEC PDP-10 at the University of Utah. This beginning gave a good indication of the diversity of host menufacturer types and operating systems which would use the ARPANET. There follows a recent- schemetic map of the ARPANET (Figure 17) showing the verious hosts. Notice that the hosts range from small PDP-11s to large IBM systems, with a scattering of very special hosts such as ILLIACEIV, running a veriety of operating systems from relatively standard ones provided by manufacturers to very special ones constructed by university researchers.

Figure 18 sharts the increase in number of hosts on the network against time.

Figure 19 provides a breakdown of the numbers and kinds of hosts on the network. It is based an information complied by the NIC for inclusion in the ARPANET directory. Same of these hosts there a single host port on an IMP. Also, each of the twenty-three TIPs in the ARPANET logically provides a host

^{*} A sampling of such aghematic maps covering the entire history of the ARPANET is provided in Appendix 1.C.

Section 1

DRAFT

Figure 17 -- ARPANET Legical Map at June 1977

function although no physical host separate from the network node is required.

Figure 18 -- Number of ARPANET Hosts Vs. Time

Figure 19 -- Summery of ARPANET Hosts

1.5 The Impact of the ARPANET

The ARPANET has served well its original function as a testbed for a new computer communications technology. More recently the ARPANET has also given goed operational service to a number of users who have some to depend on it for their computer communications service. However, a part of the original program plan for the ARPANET was technology transfer. The original program plan stated, "The transfer of interestive samputer network technology will occur in three forms; (1) Dissemination a f techniques and experimental results through the open selentific and technical literature. (2) Through the common carriers or other commercial organizations concerned with deta transfer and dissemination, and (3) Through the military command and control century for which the national Military Command System Support Center in the Pentagga serves as the focal point." The ARPANETic greatest suscess has perhaps been in this cree of technology trensfer.

Being an unclassified effort, implemented for the most part by individuals with an academic or research leaning, there have naturally been numerous papers written on the ARPANET. Two key sets of papers written relatively caply in the ARPANET

development were a set of five presented in a session at the AFIPS 1978 Spring Joint Computer Cenference and emother set of Five presented in a mession at the AFIPS 1972 Spring Joint Computer Conference. Both these sessions were organized by IPT and the two sets of five papers were specially bound together under ARPAeprovided govers and distributed widely. already listed these and a number of other papers in the bibliographies provided elsewhere in this history. In 1975 IPT commissioned decker and Hayes, Inc., of Les Angeles, California, to prepare a bibliography of publications related to the ARPANET (Saleetes Sibilospassy and Index'th Subilentions ... about . TARRANEI, Becker and Hayes, Inc., February 1976, 185p.). This bibliography lists 561 relevent documents and includes a subject index. The document is evellable from NTIS under accession AD-A026900. A complete set of the papers in the bibliography was also gotlested on migrofishe. There was also a large number of informal warking papers distributed amoung the various groups and individuals working on the ARPANET development. Many of these are also covered in the Secker and Hayes bibliography. The Netional Sureeu of Standards has also senstructed a bibliography computer communications which includes hundreds of ARPANET-polated publications.

In addition to the many papers and presentations that have been given on the ARPANET, there have been demonstrations of the technology. Of course many of demonstrations have been informal or relatively low key or relatively short, but in several instances the demonstrations have been truly major productions. Recently for example, IPT, with help from several members of the ARPANET community, took & demonstration team to Europe where a series of demonstrations was given for members of the NATO staff over a two-week periods Earlier, for the First International Conference on Computer Communications held in Washington, D.C., in 1972, 50 Kb phone Times were leased from existing network sites to the sonference mite at the Washington Hilton Hotel, and an ARPANET TIP was agt up in a demonstration room in the hotel for the duration of the conference. This event was the idea of Lawrence Roberts and was enchestrated by Robert Kahn. Oceans of members of the ARPANET dommunity were involved. Manufacturers of all manner of computer terminals were invited to semmes their terminals to the demonstration TIP. Throughout the conference hours, each day of the conference, there were individuals evailable to demonstrate use of programs on ARPANET heats from the manufacturer-provided terminals. A relatively thick booklet was

written, and many copies made available at the conference. by means of which visitors to the demonstration area could follow "do it vourself" directions to use programs on many network A twenty or thirty minute motion picture about the ARPANET and the promise of computer communications was produced end shown at the conference. Coming at a time when the TIP had not been availble for a very long time, when only a limited terminals had been tried with the ARPANET, and when number of many hosts had completed the initial implementation of the nedeseary host software but few had had it running for very long, the ICCC demonstration provided an important stimulus for the ARPANET community to pull together and get the network in true specational shape. The demonstration itself was a spectagular euccase; with everything working amazingly well, many visitors remarked that the ARPANET technology #reelly is real# and parried this impression beek home with them. The essurance with which Roberts premised the demonstration and the reutine way in which spoke of it while it was happening no doubt enhanced the impression taken home by the visiters, and belied the crash efforts and feelings of penis of the members of the ARPANET sommunity who were called upon to execute the demonstration,

There has been good success in transferring the ARPANET technology to other parts of the Department of Defense. The Defenge Communications Agency procured two small networks; essentially identical to the ARPANET in function, for the purpose of gaining experience with the ARPANET technology. Called the PWIN and EDCN networks, the first of these was used in the WWMCCS effort and the segond was used in connection with the successor to AUTODIN I. NSA also procured two smaller networks essentially identical to the ARPANET. Called COINS and PLATFORM, the first is used to connect the computers of a number of intelligence agencies and the second is used for internal NSA computer communications.

Where direct gopies of the ARPANET have not been progured, the ARPANET technology has nonetheless effected the sharesteristics of new DoD networks being built. The new SATIN IV, the Strategie Air Commend network presently under construction, has a strong component of packet switching in its makeup. The new DoD common user network, AUTODIN II, being sonstructed by DCA, is explicitly a second generation ARPANET.

Dutside the U.S. military, the semmeralal world has begun to use the ARPANET technology or veristions on it. Several

companies have filed with the F.C.C. or already been ligenced to offer communications services based more or less directly on the pecket-switching technology developed in ARPANET. Among these are Telenet Communications Corporation (for which BBN arranged the financing and Lawrenge Roberts was President and is now Chairman of the Soard), Tymnet, Graphnes, ITET, and ATET. Because of its access to substantial ARPANET expertise, of the severel common gerriers Telenet has used the ARPANET technology most directly. Telenet now serves about eighty dities in the gontinental U.S. and has made arrangements to connect to several foreign networks and serve several foreign cities, initially developed in perallel to ARPANET as a means of meking the time-sharing services of Tymshare available to a wide geographic area, and used a substantially different although related communications tehenology; however, a new version of Tymnet being built uses techniques closer to those developed for ARPANET, Graphnet, 1751, and ATST all have announded their intention to provide public packeteswitching services.

A number of U.S. compenies have also produced or are producing private corporate networks utilizing many of the techniques developed for ARPANET. For instance, it was recently announced that Citibank of New York City has constructed (by

contract to BBN) a private network very similar to the ARPANET, an increasing number of domercial RFPs call for packetsewitching or for functions which can only be provided using packet switching. A number of dompanies have taken advantage of the fact that the ARPANET technology is in the public domain to obtain the listings of the ARPANET software. There has even been a "pop" reaction to the technology. The December 7, 1972, issue of Rollingistone features an article which "approves" the popular usefulness of the ARPANET and other ARPANED developed technology, and more recently work has begun on networks of personal domputers.

While technology transfer to foreign institutions is farther from the intent of the ARPA charter, the wide foreign acceptance and use of the technology confirms the fundamental correctness and importance of the technology ARPA has developed. Several nations? PITS (the foreign national Postal, Telephone, and Telegraph authorities) have made a commitment to the development of packet-switching networks, there have been several foreign research networks, and several intermetional networks are being developed or are under consideration. There follows a list of

some of these networks:*

CIGALE == an operational network developed by a French government research agency.

RCP -- built by the French PTT and operational as a testbed.

TRANSPAC as under construction by the French PTT and scheduled to become operational for public use in 1978.

EPSS -- an experimental packeteswitching service built by the UK PTT which became operational in 1976.

CINE -- a pagket-suitching service operated by the Spanish PIT.

Detepag == a peqket=switching network built by the Trans=Canade Telephone System which is in the serily phases of operation.

JIPNET -- prestically a copy of the ARPANET built in Japan,

^{*} More detail on this list may be found in "Flenned New Publis Data Networks", P.T. Kirstein, <u>Camputar</u>, <u>Networks</u>, Volume 1, Number 2, September 1976, pp. 79=94; Kirstein is himself a member of the ARPANET community and was instrumental in arranging the installation of the London node of the ARPANET.

EIN == the European Information Network, built jointly by the U_aK_a , Switzerland, France, and Italy and strongly influenced by CIGALE.

EURONET -- a network under discussion by the European Economic Community, initially to use the EIN technology.

No doubt there are other networks beeides those mentioned above in the planning phase or under construction. With so networks coming into being, technology exchange and standards become important issues. From the time of the 1972 ICCC, representatives of various countries and institutions interested in computer networks met informativ to discuss their experience and to consider passible standards. In 1977 the International Network Working Group (INWG) was tormed. Modeled on the ARPANET Network Working Group, ARPA IPT essentially handepicked Vinton Cerf to be INWGfs first chairman and offered the services of the the ARPANET NIC to gooddinate and distribute INWG working notes. Later INHG begame associated with the International Federation for Information Processing, ARPA cut back its NIC support, and ARPA's influence in INWG dwindled. From its beginning, INWG was a forum of which techniques other than those used in the ARPANET were considered; nametheless, the ARPANET for a long time

remained the one big, existing network against which new ideas were compared.

The international packet—switching standardization effort has been empedicity effective. With the urging of Datepad, Transpad, Telenet, and dUKTs, GCITT, the international dommunications standards organization at which all of the worlds dommunidations authorities are represented, has with remarkable speed adopted standards for sommesting hosts to packet—switching networks and packet—switching networks to each other. Known as X.25 and X.92, these standards elegally address, for purposes of international communication, issues which were first seen to be of importance in the ARPANET development, perhaps where the ARPANET was seen to be deficient.

Even now the development sequence begun with the ARPANET is fer from complete. ARPA is currently pushing onward in such related areas as use of ground redic trensmission media in the development of a mobile masket-switching network, integration of panketeswitching technology with setellite broadcast and multiwaccess technology, development of security machanisms adequate for use with a packst-switching network, interwnetwork 01 connection. Improved hostatochost protocols. 450

pagket-switching technology in intra-field and fleet-to-shore communication; peaket transmission of speech, and so forth. Other institutions, metably Xerex, have developed a high-speed multi-eccess bus using packet-switching technology; this effort is under the direction of Dr. Robert Metgaife, serlier a member of the ARPANET community.

The mention of Metealfe calls to mind one of the most effective routes for ARPANET technology transfer as by the movement of people, Metealfe is but one of hundreds of individuals who have worked on the ARPANET or used one of its many sites, later moved to an institution removed from direct contact with the ARPANET, and helped convert the new institution to use of the ARPANET technology.

A final thought on the impact of the ARPANET, While the ARPANET continues to provide useful service, and while ARPA is pursuing a number of interesting advanced programs, the ARPANET is now length que of the mainstream of packet-switching network development. There are second-generation networks being built by the U,S. government, sommercial industry, and many foreign institutions. It is with these other, newer networks that new protocols are being developed, new ideas are being tried, and

ORAFT History

Section 1

there are large groups of eqtive researchers. The ARPANET has been passed by end is no longer the senter of the state of the art.

1.6 Maturity and Handever to DCA

From the very beginning it was assumed that ARPA would eventually give up operation of the network. The program plan for the network project written in the middle of 1968 noted that once the experiment of building the network was template, a dommon corrier dould be requested to take over the management of the network of IMPs and to provide "digital message service directly to the individual users on a terriff basish which would permit ARPA to terminate its system responsibility. It was hoped that this transfer dould take place three or four years after the beginning of the network development. ARPA still had in mind attempting to transfer the network to private industry in 1972 when the Office of Telecommunications Programing wrate a letter to ARPA strongly unging them to divest the network to private industry.

For help in figuring out how the trensfer of the network to a common gerrier should take place, IPT contracted with Paul Baran and his golleagues at Cabledata Associates, Inc. of Palo Alto, California, to study the leave. This study took place between April 1973 and January 1974, and the finel result of the study was to serve as the basis for an RFF for a common carrier

to take over the metwork. A douple of lengthy interim meports and a messive final report resulted from the study ("ARPANET Management Study", Final Technical Report, Cabledate Associates, Inc., January 14, 1974).

By the time the Cabledata study was done, the idea of spinning off the network to a private dompany had just its appeal. Possibly a government auditor had something to do with turning things toward internal dovernment divestiture rather than private when he worried out joud about the prepriety of turning over a government-developed network to a sommon carrier which would make a profit selling network service back to the government. In any case, after a good bit of searching for a proper home for the network, it appeared to ARPA that DCA might be the place.

At this point DTAACS stepped in and seted as metchmaker between DCA and ARPA. It is noteworthy that at the time, Dr. Rechtin was head of DTAACS, the same Dr. Rechtin who as Director of ARPA had signed the original program plan for the ARPANET, and who had become a real believer in the network.

The general guidelings of transfer of the network from ARFA to DCA were worked out by Col. Russell, then Deputy Director of

IPT, and Dr. Estil Heversten, then Deputy Technical Director of DCA. Their negotiations resulted in a memorandum of agreement which was signed in the first days of March 1975 by Lt. Gen. Lee Paschall, Director, Defense Communications Agency, and Dr. George Heilmeter, then Agting Director, Defense Advanced Resserch Projects Agency.

The memorandum of agreement called for management of to ARPANET to be transferred to DCA es of July 1, 1975, with a sixementh phase-ever period from July 1 until December 31, during which ARPA would continue to help DCA with ARPANET management while DCA acclimated itself to the Job. The memorandum also salled for a detailed transition plan to be written, which Stephen Walker of IPT succeeded in completing (with essistance from Science Applications Inc.) by June 1975.

The menagement of the network was officially transferred to DCA on July 1, 1975, However, speaking resistationally, Stephen Welker of IPT continued as menager for the network for several more months, with help from Mr. Robert Brownfield of DCA. As Brownfield learned the Job, he took mere and more responsibility for menagement of the network until by January 1975 he was carrying out the developeday menagement of the network with help

from Walker. Walker then datached himself from day to day involvement in the network and continued to serve only as ARPA representative recerding ARPANET affairs. Along with the transfer of ARPANET management from ARPA to DGA, the techniquifunctions that had been being performed at RML were also transferred to DGA and the progurement functions were transferred from RML to DECCO, the progurement agency with which DGA was most used to doing business.

There are several aspects of the memorandum of agreement and the transition pien which are worthy of mention here. network was to be an operational DoD facility, to be used solely for government business. The condept of MARPANET sponsore" was invented with sponsors being these users or collections of users (e.g., ARPA, NSS) who originally owned ARPANET equipment before management of it was turned over to DCA. Ownership of the equipment was to remain with the appnoors, DCA was to finance the operation and maintenance of the network through use of the DCA managed Communications Industrial Funds which would recover its costs by a propreted allegation to sponsors based on the equipment used by the appnears, DCA was to contract initially and SRI to perform the ARPANET operations and with BBN maintenence and NIC functions, and with NAC initially

topological consulting was needed; it was eleepty implied that DCA could retain other contractors to perform these functions eventually and certainly after the first years. DCA was to operate the network for a period of three years and thereefter if necessary until equivalent service sould be provided.

The transfer of the network to DCA ends the period govered by this history; however, it may be of interest to the reader that the network has new been under DCA management for over two years and gentinues to serve its users well. Plans are being made regarding the disposition of the ARPANET once the new common user network new being built by DCA is finished, but the ARPANET will remain in operation for at least another souple of years, and possibly for many more.

- 1.A BBN ARPANET Bibliography
- 1.A.1 Technical Reports
- AD682905 BBN Report No. 1763, "Initial Design for Interface Message Processors for the ARPA Computer Network", Jenuary 1969,
- #ADAG19160 BBN Report No. 1822, #Specifications for the Interconnection of a Host and an IMP*, revision of January 1976.
- AD730725 BBN Report No. 2161, "A Study of the ARPA Network Design and Performance", Kahn and Crowther, August 1971.
- *ADAU1439\$ BON Report No. 2183, "Terminal Interface Message Processor = User*s Guide," updated version of August 1975.
- *ADAGG2461 BBN Report No. 2164, "TIP Herdwere Henuel," revision of November 1974.
- *AD776995 BBN Report No. 2277, *Specifications for the Intertannection of Terminals and the Terminal IMP,* Retthera, revision of June 1973,
 - BBN Report No. 2491, "Throughput in the ARPA Network or Analysis and Measurement," McQuillen, Jenuary 1973 (text also contained in GTR No. 16, Jenuary 1973).
- ADA025356 BBN Technical Information Report (TIR) No. 89, "The Interface Message Progessor Program," updated version of June 1976 (Version 3231).
- ADAD33351 BBN TIR No. 98, "The Network Control Center Program," updated version of November 1976 (Version 131).
- ADAG31617 SSN TIR No. #1, "The Terminal Interface Message Processor Program," updated Version of August 1976 (Version 486).
- AD781467 BBN Report No. 2831, "Adaptive Routing Algorithms for Distributed Computer Networks," McQuillan, May 1974,

BBN TIR No. 93, The Remote Jeb Entry MiniwHost, R August 1974 (text also conteined in QTR No. 6, July 1974).

BBN Report No. 2891; "A Proposed Experiment in Peaket Broadcast Satellite Communications;" Rettberg and Walden, September 1974.

BBN Report No. 2915, *Network Design Issues,* November 1974 (text also genteined in QTR Ne. 7, October 1974),

ADAD22040 BSN Report No. 2999, #Fluribus Dogument 11 Overview, # May 1975.

BBN Report No. 2938, *Fluribus Desument 2: System Handbook, * January 1975.

BBN Report No. 3000, "Pluribus Document 31 Configurator," in production.

BBN Report No. 3001, "Fluribus Dosument 4: Basis Boftware," Desember 1975.

BBN Report No. 2931, *Pluribus Document 5: Advanced Software, * April 1975.

BBN Report No. 3882, "Fluribus Document by Functional Specifications," February 1976,

MBN Report No. 3684, "Pluribus Document '7; Haintenance," September 1976,

SBN Report No. 3856, "The Atlantic Satellite Pecket Broadcat and Getsway Exmeriments," R. Sinder, A. Rettberg, and D. Walden, April 1975.

ADAG1834: BBN Report No. 3126, "A Multiprocessor Design," Wa Berker, September 1975.

1.A.2 Published Paperse

W. Teitleman and R.E. Kahn, "A Network Simulation and Display Program," Progeedings of the Third Annual Princeton Conference on Information Sciences and Systems, Harch 1969.

F.E. Heart, R.E. Kahn, S.M. Drastein, W.R. Crewther, and D.C. Welden, #The Interface Message Processor for the ARPA Computer Network, # APIPS Conference Proceedings 36, June 1978, pp. 551=567; also in Advances in Computer Communications, W.W. Chu (ed.), Artech House Inc., 1974, pp. 388-316; also in Computer Communications, P.E. Green and R.W. Lucky (eds.), IEEE Press, 1975, pp. 375-391; also in Computer Networking, R.P. Blanc and I.W. Cotton (eds.), IEEE Press, 1976, pp. 66-76.

FiE: Heart and BiM: Ornstein, "Baftwere and Louis Design Interaction in Computer Networks;" International Computer State of the Art Report No. 6: Computer Networks; Infotech Information Ltd., Meidenhead, Berkshipe, England, pp. 423-462.

 $R_{\pm}E_{\pm}$ Kahn, "Terminal Aggess to the ARPA Computer Network," in Courant Computer Science Symposium 3 we Computer Networks, R_{\pm} Rustin (ed.), Prentice=Hall, Englewood Cliffs, NaJa, 1972, ph. 147=166.

R.E. Kahn and W.R. Crowther, "Flew Control in a Resource Sharing Computer Network," Presentings of the Second ACM/IEEE Symposium on Problems in the Optimization of Date Communications Systems, Palo Alto, California, October 1971, pp. 188-116; elso in IEEE Transactions on Communications, Vol. COM-28, No. 3, Part II, June 1972, pp. 539-546; elso in Advances in Computer Communications, W.W. Chu (ed.), Artoch House Inc., 1974, pp. 539-837; elso in Computer Networking, R.P. Blane and I.W. Cotton (eds.), IEEE Press, 1974, pp. 117-125.

D.C. Helden, "A System for Interpretess Communication in a Resource-Sharing Computer Network," Communications of the ACM, Vol. 15, No. 4, April 1972, pp. 2214236; also in Advances in Computer Communications, W.W. Chu (ed.), Artoch House Inc., 1974, pp. 3404344.

^{*} Ordered according to date of writing.

R₁H₂ Thomas and D₄A₂ Hendersen, *MgRD\$8 we A MultiwComputer Programming System, * AFIPS Conference Proceedings, Yol, 40, June 1972, pp. 261+293; also in Computer Networking, R₂P₆ Bland and I₂W₆ Cotton (eds.); IEEE Press, 1976, pp. 246=258,

S.M. Ornstein, F.E. Heart, W.R. Crowther, S.B. Russell, H.K. Rising, and A. Mishel, "The Terminal IMP for the ARPA Computer Network," AFIPS Conference Proceedings 40, June 1972, pp. 243-254; also in Advances in Computer Communications, W.W. Chu (ed.), Artesh House Inc., 1974, pp. 317-326; also in Computer Communications, P.E. Green and R.W. Lucky (eds.), IEEE Press, 1975, pp. 354-365.

He Franks, R.E. Kahn, and L. Kieinrocks, "Computer Communications Network Design so Experience with Theory and Fractice," AFIPS Conference Protectings, Vol. 45, June 1972, pp. 255-270; also in Networks, Vol. 2, Ne. 2, 1972, pp. 135-166; also in Advances in Computer Communications, W.W. Chu (ed.), Artech House Inc., 1974, pp. 254-269.

A.A. McKenzie, B.P. Cegeii, J.M. McGuilian, and M.J. Threps, "The Network Centrel Center for the ARPA Network," Precedings of the First International Conference on Computer Communication, Hashington, D.C., October 1972, pp. 185-191; also in Computer Networking, R.R. Bland and I.W. Cotton (eds.), IEEE Press, 1976, pp. 319-325.

R.E. Kehn, "Resource-Sharing Computer Communications Networks," Proceedings of the IEEE, Vol. 68, No. 11, Nevember 1972, pp. 1397-1407; also in Advenses in Computer Communications, Walk, Chu (ad.), Artech House Inc., 1974, pp. 208-218; also in Computer Communications, P.E. Green and R.M. Lucky (eds.), IEEE Press, 1975, pp. 537-547.

J.M. McQuillen, W.R. Crewther, 8.P. Cesell, D.C. Walden, and F_4E_4 Heart, "Improvements in the Design and Performance of the ARPA Network," AFIPS Conference Proceedings 41, Desember 1972, pp. 741-754,

enternamentum eNetwork Analysis Corporation, esuniversity of California at Los Angeles,

WaRs Growther, RaDs Rettberg, D.C. Welden, Same Ornstein, and F.E. Heart, "A System for Broadcast Communications Reservations ALOMA," Proceedings of the Sixth Hawaii International Conference on System Sciences, January 1973, pp. 371-374.

D.C. Welden, "HostotemHost Protocels," International Computer State of the Art Report No. 24: Network Systems and Software, Infotech, Maidenhaed, England, pp. 287=316.

N.W. Mimne, B.P. Cosell, D.C. Melden, B.C. Butterfield, and J.B. Levin, "Terminal Access to the ARPA Network so Experience and Improvements," Proceedings of the Seventh Annual IEEE Computer Society International Conference, San Francisco, Colifornia, February 1972, pp. 30-43; also in Computer Networking, R.P. Bland and I.W. Cotton (edg.), IEEE Frence, 1976, pp. 287-291.

E.W. Wolf, "An Advanced Computer Communication Network," AIAA Computer Network Systems Conference, April 1973, AIAA Paper No. 73-414.

P.E. Heart, S.M. Ornstein, W.R. Crewther, and W.B. Barker, "A New Minisomputer/Multiprocessor for the ARPA Network," AFIPS Conference Precedings 42, June 1973, pp. 329-537; also in Advances in Computer Communications, W.W. Chu (ed.), Artach House Inc., 1974, pp. 329-337; also in Computer Communication Networks, R.L. Brimsdele and F.F. Kuo (eds.), Proceedings of the NATO Advanced Study Institute of September 1973, Sussex, England, published by Neerdheff Internations; Publishing, Leyden, The Natherlands, 1975, pp. 189-168; also in Computer Communications, P.E. Green and R.W. Lugky (eds.), IEEE Press, 1975, pp. 366-374,

R.H. Thomas, "A Resource Sharing Executive for the ARPANET," AFIPS Conference Presentings 42, June 1973, pp. 153-163; also in Advances in Computer Communications, M.W. Chu (ed.), Artech House Inc., 1974, pp. 359-367.

F.E. Heart, The ARPA Network, In Computer Communication Networks, R.L. Grimsdelp and F.F. Kud (eds.), Proceedings of the NATO Advanced Study Institute of September 1973; Sustan, England, published by Noordheff International Publishing, Leyden, The Netherlands, 1975, pp. 19835.

 W_2R_a Crowther, J_aM_a MeQuillan, and D_aC_a Walden, "Reliability Issues in the ARPA Network," Progeedings of the ACM/IEEE Third Data Communications Symposium, Nevember 1973, pp. 159e160; also in Computer Networking, R_aP_a Blanc and I_aW_a Cotton (ads.), IEEE Press, 1976, pp. 142e143.

J.M. McQuillen, "Decign Considerations for Rauting Algorithms in Computer Networks," Proceedings of the Seventh Annual Hawaii International Conference on System Sciences, Honolulu, Hawaii, January 1974, pp. 22m24; also in Computer Networking, R.P. Sland and I.W. Cotton (eds.), IEEE Press, 1976, pp. 1884119.

S.C. Butterffeld, R.D. Rettberg, and D.C. Walden, "The Satellite IMP for the ARPA Network," Proceedings of the Seventh Annual Hawaii International Conference on System Sciences, Honolulus Hawaii, January 1974, Computer Nets Supplement, pp. 78-75.

S.M. Drostein, W.D. Berker, R.D. Breesier, W.R. Crowther, F.E. Heart, M.F. Kreiey. A. Mithel, and M.J. Thrope, "The BBN Multiprocessor," Presendings of the Seventh Annual Hawaii International Conference on System Selences, Honolulu, Hawaii, January 1974, Computer Nets Supplement, pp. 92095.

A.A. McKenzie, "Some Computer Network Interconnection Issues," AFIPS Conference Precedings 43, May 1974, pg. 857m859,

 E_aA_a Akkeyunium, A_aJ_a Sermsteinm, and R_aE_a Schentz, "Interprocess Communication Facilities for Network Operating Systems," IEEE Computer, Val. 7, No. 6, June 1974, pp. 46=55.

F,E, Heart, "Implications of the ComputersCommunication Pertnership," in Preprints of Papers, MEDINFO 74; The First world Conference on Medical Informatics, Stephelm, Sweden, August 1974, pp. 21-27; precedings to be published by NertheHelland, The Netherlands.

adtate University of New York at Stony Brook.

Fig. Heart, "Networks and the Life Sciences; the ARPA Network and Telenet," Federation Proceedings, Federation of American Societies for Experimental Siclogy (FASES), Vol. 33, No. 12, December 1974, pp. 2399=2402; also in Computers in Life Science Research, M. Siler and D.A.B. Lindberg (eds), FASES and Plenum Press, 1975, pp. 284=215.

J.M. McGullen, "The Evolution of Message Processing Techniques in the ARPA Network," International Computer State of the Art Report No. 241 Network Systems and Seftware, Infetech, Maidenhead, England, pp. 541-576.

Rame Thomas, "JSYS Traps = A TENEX Machenism for Encapsulation of User Precessor," AFIRS Conference Proceedings 44, May 1975, pp. 351=360.

S,M. Ornstein, W.R. Crowther, M.F. Kreley, R.D. Bressler, A. Michel, and F.E. Heprt, "Pluribus we A Reliquie Multiprocessor," AFIPS Conference Proceedings 44, May 1975, pp. 551-559.

F.E. Heart, S.M. Ornstein, W.R. Crowther, W.S. Berker, M.F. Kraley, R.D. Bressler, and A. Michel, "The Pluribus Multiprocessor Systems" in Multiprocessor Systems: Infotech State of the Art Report, Infotech International Ltd., Maidenhead, Berkshire, England, 1976, pp. 307=336.

WaRs Crewther, F.E. Heart, A.A. McKenzie, J.M. McQuillen, and D.C. Walden, "Issues in Packet=Switching Network Design," AFIPS Conference Precedings 44, May 1975, pp. 161=175; else in Computer Networking, R.P. Blane and I.W. Cotton (eds.), IEEE Press, 1976, pp. 182=196.

J. Burchfiel, R. Temlinson, and M. Beeter, "Functions and Structure of a Facket Radio Station," AFIPS Conference Proceedings 44, May 1975, pp. 245-251.

R.S. Temiinson, "Seiesting Sequence Numbers," Proceedings of the ACM SIGCOMMUSIGOPS Interface Workshop on Interprocess Communications, Merch 1975, pp. 11-23.

Jama Maguijian and D.C. Welden, "Some Considerations for a High Performance Message-Based Interprocess Communication System," Proceedings of the ACM SIGCOMM#SIGOPS Interface Workshop on Interprocess Communications, March 1975, pp. 77-86.

Sing Ornstein and D.C. Maiden, "The Evolution of a High Performance Meduler Meeket=Switch," Conference Record of the 1975 International Centerence on Communications, June 1975, Vol. 1, pp. 6=17 to 6=21.

Rabinder, MA Dynamia Packet=Switching System for Setallite Broadcast Channels, Conference Record of the 1975 International Conference on Communications, June 1975, Vol. III, pp. 41w1 to 41w5.

R,D, Bressier, M,F, Xreisy, and A, Mishel, MPluribust a Multiprosessor for Communications Networks, M Paurisonth Annual ACM/NSE Technical Symposium we Computing in the midw78fs; an Assessment, June 1975, pp. 13w19.

D.C. Welden, "Experiences in Building, Operating, and Using the ARPA Network," Proceedings of the 2nd USA-Japan Computer Conference, Tokyo, Japan, August 1975, pp. 453-458.

RaDa Rettberg and DaCa Walden, "Gateway Design for Computer Network Interconnection," int Gommunications Networks: Proceedings of the European Communications Conference on Communications Networks, London, England, September 1975, published by Online Conferences Limited, Umbridge, England, pp. 113-126.

M.F. Kraley, "The Pluribus Multiprocessor," Digest of the 1975 International Symposium on PaultoTelerant Computing, Paris, France, June 18-25, 1975, p. 251 (abstract only).

 B_aP_a Ceasily J_cM_a MeQuillan, and D_aC_a Majden, "Techniques for Detecting and Preventing Multiprogramming Bugs," int Minisomputer Software, J_aR_a Sell and C_cG_a Sell (eds,), Northeholland Publishing Co., 1976, pp. 381=388.

A,A, McKengie, "The ARPA Network Control Center," Proceedings of the Fourth ACM Data Communications Symposium, Quebec City, Canada, Ostober 1975, pp. 5=1 to 5=6.

R. Binder, J.M. McGuillen, and R.D. Rettberg, "The Impact of Multi-Access Extellites on Packet Switching Networks," EASCON 775 Record, IEEE Electronies and Aerospece Systems Conference, Washington, D.C., October 1975, pp. 63-4 to 63-7.

 B_1P_2 Cosell, P_1R_4 Johnson, J_4M_4 Melmen, R_1E_2 Schentz, J_4 Sussmen, R_2H_4 Thomes, and D_4C_4 Melden, "An Operational System for Computer Resource Sharing," Precedings of the Fifth Symposium on Operating System Principles, Austin, Texas, November 1975, pp. 75-81.

RiD, Bressier and R.H. Thomas, "Design Issues in Distributed Computing," in Distributed Systems: Infotesh State of the Art Report: Infotesh International Ltd., Maidenhead, Berkshire, England, 1976, pp. 215-222.

V. Cerfe, A. McKengle, R. Scentlebury**, and H. Zimmermenn***, **Proposal fer an International End to End Protosol, ** ACM Computer Communication Review, Vol. 6, No. 1, January 1976, pp. 63*69.

WaFe Mann, 8.M. Ornstein, and MaFe Kraiwy, TA NetworksOriented Multiprocessor Frontagnd Handling Many Heats and Hundreds of Terminals, AFIPS Conference Proceedings 45, June 1976, pp. 533-548.

P_iJ₂ Santoe, "Software Instrumentation for Maintainability of Distributed Computer Networks," Fifteenth Annual ACM/NBS Technical Symposium-Directions and Challenges, June 1976, pp₂ 143-146,

R. Binder, "A Reservation-TDMA Protess! for the Atlantic Packet Satellite Experiment," Workshop on Packet-Switching and Communication Satellites, Institut de Recherche d'Informatique et d'Automatique, Rennes, France, June 1976 (Abstract only).

^{*}Defense Advanced Research Projects Agendy
**Institut de Recharche d'Informatique et d'Autematique (France)
***National Physical Laboratories (U,K,)

 R_{\pm} Binder, Nateling May, $R_{\pm}D_{\pm}$ Rettberg, $D_{\pm}G_{\pm}$ Walden, and $R_{\pm}G_{\pm}$ Walsaler, "Current Status of the Satellite Interface Message Processor," Warkshop on Pasket=Switching and Communication Satellites, $I_{\pm}R_{\pm}I_{\pm}A_{\pm}$, Rennes, France, June 1976 (Abstract only),

J'M, McGuillen, MA Status Report on the ARPANET, M Australian Computer Society/International Federation for Information Processing Jaint International Symposium on Data Communications Technology and Practice, August 1976, pp. 13,1013,7.

J.M. McQuillan, "Strategies for Implementation of Multi-Host Computer Networks," ACS/IFIP Joint International Symposium on Date Communications Technology and Practice, August 1976, pp. 26,1=26,6; also in Computer Communication Review, Vol. 6, No. 4, October 1976, pp. 19=24.

Frequence and D.C. Welder, "Communications Applications of the Pluribus Computer," Conference Record of the 1976 IEEE National Telecommunications Conference, Dallac, Texas, 29 November 1 December 1976, pp. 7,1=1 to 7,1=5.

As Evens and G_0R_0 Morgans "The Tinman and Communications Applications," presented at the Workshop on the Design and Implementation of Pregramming Languages; Cornell University, October 1, 1976,

 R_s Levine, J_sM_s McQuillan, and R_s Schentz, "Distributed Systems," a section of "New Directions for Operating Systems: A Workshop Report," J_sC_s Browns, Operating Systems Review, Vel. 11, No. 1, January 1977, pp. 14=19.

 J_*M_4 McGuillan and D_*G_8 Walden, "Issues in Packet Switching Network Design and the ARPANET Design Designap," submitted to Computer Networks.

J. Davidsen, W. Hathaways, N. Mimno, J. Postelas, R. Thomas, and D. Walden, "The ARPANET TELNET Protogols Its Purpose, Principles, Implementation, and Impact on Hest Operating System Design," submitted for presentation at the ACM/IKEE Fifth Data

^{*************}

^{*}Carnegie=Hellon University

anasa Ames Research Conter

^{**}USC Information Sciences Institute

Communications Symposium, Snowbird, Utah, September 27-29, 1977.

J. Jacobanne, L. Leenne, A. Viterbinne, R. Binder, R. Bressler, N. Hau, and R. Heissler, "CPODA -- a Demand Assignment Protocol for SATNET," submitted for presentation of the ACM/IZEE Pifth Date Communications Symposium, Snewbird, Utah, September 27-27, 1977.

***Linkabit Corporation

1.B Host Protecol Chronology

Network Working Group formed, First RFC published, 2. Crocker "Hest Boftware," RFC 1. NIC 4687, 7-Apr=69. May 1969 First heat level Protocel Specified.
G. Deleche "Hoat Seftwere," RFC 9, NIC 4695, 1=May=69. February 1978 Second host level protocol proposed. 5. Creeker "New HOSTeHOST Protocol," RFC 35, NIC 4735, 12-Fab=78. June 1970 Versian 2 of second hest level protects proposed. 8. Crocker, "An Official Protocol Proffering." RFC 54. NIC 4756, 18-Jun-79. August 1978

Host level protogel published as official specification. & Crocker, "Official Host to Most Protocol," Document 1, NIC 7149, 3-Aug-70.

The Message Switching host level protocol is proposed, D. Walden AA System for Interpresent Communication in a Resource Sharing Computer Network, # RFC 62, NIC 4962, 3-Aug=78, Alse published in Communications of the ACM, 15(4);221=230, April 1972.

Network Graphies Meeting

Where , Who 777 Initial Connection Protocol proposed.

8. Creeker #3rd Level Ideas and other Nelse, # RFC 66, NIC 5489, 26-Aug-70.

November 1978

Network Werking Group Meeting at FJCC in Houston, J. Pestel "Network Meeting Report," RFC 77, NIC 5684.

ZdeNeve7e, E. Mever "Network Meeting Netes," RFC 82, NIC 5619, 9-0-4-79.

Initial network use at UCSS and RAND.

E. Harmiem "NCP Status Report: UC\$B/RAND, " RFC 78, NIC 5199, undated.

Desember 1971

Data Reconfiguration Service first suggested,

E. Hersiam "Protocols and Deta Formets," RFC 80, NIC 5608, 1-Dec-70.

January 1971

```
Graphies Protocol first proposed.
      8. Gracker "Proposal for a Network Standard Format for a
      Data Stream to Centrol Graphics Display," RFC 86, NIC 5631,
      SeJane71.
   Remote Jeb Entry protocol first proposed.
      R. Braden "NETRJS = A Third Lave! Protoco! for Remote Job
      Entry, " RFC 48, NIC 3648, 130Jano71,
   Initial use of the network at MIT and Hervard,
      R. Meteelfe "Some Historia Momenes in Natworking." RFC 89.
      NIC 5697, 19-Jan-71.
February 1971
   Teines protecti first proposed,
      J. Melvin MA Firt Cut at a Proposed Talmet Protocol." RFC
      97. NIC $748. 15mFeb=71.
      E. Mayer "Lagger Protocol Proposed," RFC 98, NIC 3744,
      11efebe71.
   Network Working Group held at University of Illinois.
      R. Watson "Notes on the Network Working Group Meeting," RFC
      181, NIC 5742, 23=Feb=71,
      R. Wateen, "Addendum to NWG Meeting Notes," RFC 188, NIC
      5887, 25-Mar-71,
   Host level protect medified,

    Crocker "Output of the Host/Host Protocol Glitch Cleaning

      Committee." RFC 102, NIG 5763, 22=Feb=71.
March 1971
   Host level protocol medified.
      R. Bressier "Output of the HestaHost Protocol Glitch
      Cleaning Committee, # RFC 187, NIC $886, 23-Mar-71.
April 1971
   File Transfer Protogol first proposed.
      A. Bhyshen "A file Transfer Protocol," RFC 114, NIC 5823,
      16-APP-71.
May 1971
   Network Working Group meeting held at SJCC in Atlantic City.
      J. Heafner, "Ninutes of the Network Working Group Mesting,"
      RFC 144, NIC 4778, 25-May=71,
   Telnet protocol specified.
      T. O'Builivan TTELNET Protocol, " RFC 158, NIC 6768,
   19=May=71.
Data Resontiguration Service specified.
      B. Anderson *Data Reconfiguration Spevise we An
      Implementation Specification, RFC 146, NIC 6780, 25-May=71.
June 1971
   File Transfer Protocol design initiated,
```

```
Bhusham, A. The Data Transfer Protocol, RFC 171, NIC 6793,
       21-JUN-71.
       Shusham, A. "The File Transfer Protocol," RFC 172, NIC 6794,
       23-JUN-71.
July 1971
   Mail Protecol first discussed.
       Watson, R. WA Meil Sex Protosol, W RFC 196, NIC 7141,
       28-JUL-71.
   Graphics meeting held by MIT.
August 1971
   Terminal IMP Protocol Implementation
      McKeng(e, A, MNCP, ICP, and Telnet: The Terminal IMP Implementation, M RFC 215, NIC 7945, 30-Aug-71.
October 1971
   Network Working Group Meeting and Programmers Workshop
      Vesse, A. *Network Working Group Meeting Schedule, * RFC 234,
      NIC 7651, 5-0CT-71,
      Postel, J. "Report of the Protocol Werkshop," RFC 295, NIC
      8335, 2-JAN-72.
November 1971
   Notwork Graphics Meeting
      Padlipaky, M. "Graphies Meeting Report," RFC 282, NIC 8164,
      S-DEC-71.
January 1972
   Protocol proposed by Graphica Committee.
      Michener, J. "Grephies Protocol - Level @ only." RFC 292,
      NIC 8302. 12-JAN-72.
April 1972
   file Transfer protocol Meeting
      Shushan, A. "Deta and File Transfer Workshop Nates," RFC
   327, NIC 9634, 17-MAR-72. Graphics protocol meeting
   Teinet Protecol appelfication published
      Postel, J. "Telnet Protocol," RFC 318, NIC 9348, 3-APR-72.
   Remote Job Entry protocol meeting
June 1972
   Remote Job Entry protecti subjished
      Helland, C. "Proposed Remote Job Enter Protocol," RFC 360,
      NIC 18682, 24-JUN-72.
August 1972
   Meil commands included in evolving File Transfer protocol
   specification
      Bhushen, A. *Comments on the File Transfer Protosol (RFC
      354), " RPC 385, NIC 11357, 18-AUG-72,
```

```
October 1972
   Demonstration of ARPANET at International Computer
   Communications Conference (ICCC) in Washington D.C.
January 1973
   The principle of negotiated options for telnet first
   enungieted.
      Cosell, 8, "TELNET Issues," RFC 435, NIC 13675, 5-JAN-75.
February 1973
   A semmen user command lenguage proposed.
      Redlineky, M. *Tenetive Propessi for a Unified User Lavel
      Protocol." RFC 451, NIC 14135, 22-FEB-73,
   Mail protogol meeting
      Kudilek, M. "Network Mail Meeting Summery, P RFC 469, NIC
      14798, 8-MAR=73.
      Shushan, A. FFTP and Network Mail System, " RFC 475, NIC
      14919, 6-MAR-73.
March 1973
   Teinet Pratecol Hesting
      Mckensie, A. "Teinet Protoco! Meeting Announcement," RFC
      461. NIC 14416, 14-PEB-73.
   File Transfer protocal Meeting
      Mekansie, A. Affie Transfer Protesti Meeting Annaungement
      and a Proposed Document, RFC 454, NIC 14333, 16mFEBe74,
May 1973
   Telnet protocol epocified, a thorugh revision of the control
   signal machanisms and the adoption of the negotiated option
   concept.
      Mekenzie, A. "Telnet Pretocol Specification," RFC 495, NIC
      15371, 1-MAY-73
   Resource Shering Executive Workshop meeting
   Usens Interest draup meeting.
      Crocker, D. Farpenet Users Interest Working Group Meeting,"
      RFC 565, NIC 20050, 6-NOV-73,
June 1973
A file access protocol is proposed.
      Day, J. *File Access Pretegal, # RFC 528, NIC 16819,
      25-JUN1-73.
July 1973
   Network Graphics Heeting
      Sundh, S. "Minutes of Network Grephies Group Meeting," RFC
      549, NIC 17795, 21-AUG-73.
October 1973
   A common text editing subsystem proposed, Padiipsky, M. "Neted: A Common Editor for the ARPA Network,"
```

RFC 569, NIC 18972, 15-0CT=73,

July 1974

A cross-network debugging program is proposed,
Meder, E. "Network Debugging Protocol," RFC443, NIC 30573,
JUL-74.

November 1974

A Host to Front End protegol preposed:

Padlipsky, M. TA Proposed Protogol for Connecting Host
Computers to ARPAWLike Networks via DirectlysConnected Front
End Processors, PRFC 647, NIG 31117, 12=NOV=74.

Desember 1974

A procedure oriented protocol is preposed.

Postel, J. **Procedure Call Protocol Documents,** RFC 674, NIC 31484, 12**DEC=74.

June 1975

Mejer extensions to the IMP/Host interface proposed Welden, D.C. *IMP/Host and Host/IMP Pretocol Change,* RFC 687, NIC 32654, 2=JUN-75.

DRAFT History

1.C Selection of ARPANET Logical Maps

We can select from the followings

December 1978 February 1971 April 1971 August 1971 December 1971 March 1972 August 1972 February 1973 August 1973 September 1973 Jenuary 1974 June 1974 November 1974 January 1975 April 1975 June 1975 July 1975 October 1975 July 1976 August 1976 October 1976 November 1976 December 1976 Merch 1977 June 1977

2. DESIGN AND IMPLEMENTATION

This chapter of our report describes the technical basis for the ARPANET, We begin with an overview of some of the main technical concepts of the ARPANET and of packateswitching. Then, we so on to treat in turn the tepological design of the system, the nature of the communications equipment within it, the subnetwork protocols, and the host level pretocols, and finally, a treatment of the performence of the overall ARPANET system.

2.1 Overview and Technical Bests

The ARPANET represents e major step forward in computer communications. As such, is introduced a number of new terms into the technical literature. We will discuss these terms and present the design point of view which underlies much of the ARPANET technology. We present this parapagive in two pertst first, a discussion of the design goals of the system, and secondly, a brist overview of the sholess made in each of the important design areas in the ARPANET.

2.1.1 Terminology

Mades. The nodes of the network are real-time computers, with I (mited storage and processing resources, which perform the basic packet-switching functions.

MERIE. The Hosts of the network are the computers, connected to nodes, which are the providers and users of the network services.

Liber. The lines of the network are some type of communications circuit of relatively high bandwidth and reasonably low error rate.

Casing the degenerate cases of the general topology we consider.

MRESSOR. The unit of data exchanged between source Host and destination Host.

Ragket. The unit of data exphanged between adjagent modes,

askaguiedansi. A piece of control information returned to a squrce to indicate successful reselpt of a packet or message. A packet acknowledgment may be returned from an adjacent node to indicate successful reselpt of a packets a message adknowledgment may be returned from the destination to the squrce to indicate successful reselpt of a message.

Store and Entward Subsettant. The mode stores a copy of a packet when it receives one, forwards it to an adjacent mode, and discards its copy only on receipt of an acknowledgment from the adjacent mode, a total storage interval of much less than a second.

Parket, Switching. The nodes forward packets from many sources to many destinations slong the same line, multiplexing the use of the line at a high rate.

Revise. Alserithm. The procedure which the modes use to determine which of the severel possible paths through the network will be taken by a pecket.

NederNade: Itansbissies, Presedures. The set of procedures governing the flow of packets between edjecent nodes.

Source Destination... Inequalization... Residuent. The set of procedures governing the flow of messages between source node and destination node.

Mestaboge Transmission Presedures. The set of procedures governing the flow of information between a Host and the node to which that Host is directly connected.

HARITHOSY. Transmission Protections. The set of procedures coverning the flow of information between the source Host and the destination Host.

Some of the more obvious differences emong packetwawitching networks can be sited briefly. The ARPANET apilts messages into packets up to 1980 bits long; some of the other networks have 2000-bit packets and he multipacket messages. Hosts connect to e single node in the ARPA Network and SITA; multiple connections are possible in Cyclades and EIN. Dynamic routing is used in the ARPANET and EIN; a different adeptive method is used in SITA; fixed routing is presently used in Cyclades. The ARPANET delivers messages to the destination Host in the same sequence as it accepts them from the source Hosts Cyclades does nots in EIN the sequence is optional. Clearly, many of the design choices made in these networks are in conflict with each other. The resolution of these conflicts is essential to the planning and building of balanced, high performance networks, particularly singe many future designs will be intended for networks which are larger, less experimental, and more complex.

We next summerise how the IMP in the ARPANET performs its functions as a message switching center and interface between host computers. Figure 2=1 shows a diagram of message flow in

Section 2 Design and Implementation

Figure 2mi Packets and Messages

the ARPANET and illustrates some of the terminology. The host sends the IMP a message with up to 8863 data bits. The source breaks this up into pagkets with up to 1008 data bits. When IMP successfully received packet is at each IMP. estpowledGment or ask is sent back to the previous IMP. When the message arrives at the destination IMP it is reasonbled, that is, the packets are combined into a message again. The message is sent to the destination host and when it has been accepted, a <u>Ready For Next Hesses</u> which we abbreviate as <u>RENM</u> is sent back to the source host. The RFNM is also a packet and it is acknowledged. Several points are worth noting. First, acks are separate transmissions, but are piggywhacked in not actually packets to cut down on overhead. Next. packets on the 1 inks between IMPs are checksummed in the modem interface hardware and the IMP employs a positive acknowledgment retransmission scheme: that is, if a packet is in error, it is not acknowledged. It is then retrenemitted until an acknowledge is received. Further, an IMP may send the several peakets of a message out on different links. Sequipe of retrensmission (out of order) of a packet on a link and transmission of mackets on elternate links, the packets of a message may applied at the destination IMP out of order and must be reassembled into the correct order for transmission into the host.

2.1.2 Design Goelst Performance and Functions

In this section we define what we believe are fundamental properties and requirements of packet-switching networks and what we believe are the fundamental priteria for measuring network performance.

2,1,2,1 Basto Issues

by alvina the properties beatn We dentral packet=switching network design. The key assumption here is that processino algorithms the packet (Poutings acknowledgment/retransmission strategies used to transmission over noisy gircuits, etc.) result in a virtuel network path between the hosts with the Pallowing characteristics!

- a. Finite, flustuating delay we A result of the basis line bandwidth, speed of light delays, queueing in the nodes, line errors, etc.
- b. Finite, flustuating bandwidth == A result of network overhead, line errors, use of the network by many sources, ets.

c. Finite packet error rate (duplicate or lost packets) ==
A result of the acknowledgment system in any
storemendeforward discipline (this is a different use of
the term Merror retek then in traditional telephony).

Duplicate packets are saused when a nede goes down after receiving a packet and forwarding it without having sent the acknowledgment. The previous node then generates a duplicate with its retransmission of the packet. Packets are lost when a node goes down after receiving a packet and acknowledging it before the successful transmission of the packet to the next node. An attempt to prevent lost and duplicate packets must fell as there is a tradeoff between minimizing duplicate packets and minimizing lost packets. If the nodes avoid duplication of packets whenever possible, more packets are lost. Conversely, if the nodes retransmit whenever packets may be lost, more packets are duplicated.

d. Disordering of packets we A property of the acknowledgment and routing elgorithms.

These four properties describe what we term the store-and-forward subnetwork.

There are also two basis problems to be solved by the source and destination in the virtual path described above:

- e. Finite storage on A property of the nodes,
- F. Differing source and destination bandwidths we Largely a property of the hosts.

[Note: the question is frequently relaid whether the source and destination nodes or the source and destination hosts should solve these problems. This question is addressed in a later section.]

The fundamental requirements for packetsswitching networks are distated by the six properties enumerated above. These requirements include:

a. Buffering we Buffering is required begause it is denotedly necessary to send multiple data units on a communications path before receiving an adknowledgment. Secause of the finite datay of the network, it may be desirable to have buffering for multiple packets in flight between source and destination in order to ingress throughput. That is, a system without adequate buffering may have unacceptably low throughput due to long delays weiting for agknowledgment between

Design and Implementation

transmissions. Buffering is else required when input traffic momentarily expects autput dapadity.

- Pipelining we The finite bendwidth of the network may b. necessitate the pipelining of each message flowing through the natwork by breaking it up into pagkets in order to degrease delay. The bandwidth of the sircuits may be low enough so that $ar{t}$ orwarding the entire message at each node in the path results in excessive dalay. By breaking the message into packets, the nodes are able to forward the first peaket of the message through the natwork ahead of the later ones. For a message of P packets and path of H hops (with small propagation delays), the minimum delay is proportional to P + H = 1 instead of $P + H_e$ where the proportionality constant is the packet length divided by the transmission rate. (See section 2.5.3 below for a derivation and more exact result.)
- c. Error Control == The node=to=node pagkst processing algorithm must exercise error control, with an acknowledgment system in order to deal with the finite packet error rate of the eirquita. It must also detect when a dirquit becomes unuseble, and when to begin to use it again. In the source-to-destination processing algorithm, the destination may need to Message exercise some controls to detect missing and duplicated messages

or portions of messages, which would appear as incorrect date to the end user. Further, ecknowledgments of message delivery or non-delivery may be useful, possibly to trigger retransmission. This mechanism in turn requires error control and retransmission itself, since the delivery reports can be lost or duplicated. The usual technique is to assign some unique number to identify each date unit and to time out unenswered units. The error correction mechanism is invoked infrequently, as it is needed only to recover from node or line failures.

- d. Sequencing -- Since packet sequences can be received out of order, the destination must use a sequence number technique of some form to deliver messages in correct order, and packets in order within messages, despite any scrembling effect that may take place while several messages are in transit. The sequencing mechanism is frequently invoked since it is needed to recover from line errors.
- e. Storage allocated at either the mender or the receiver.

f: Flow Control == The different sounce and destination data rates may necessitate implicit or explicit flow control rules to prevent the network from becoming gongested when the destination is slower than the sounce. These rules can be tied to the sequencing mechanism, with no more messages (packets) accepted after a certain number, or tied to the storage allocation technique, with no more messages (packets) accepted until a sertain amount of storage is free, or the rules can be independent of these features.

In satisfying the above six requirements, the algorithm often exercises contention resolution rules to allegate resources among several users. The twin problems of any such facility are:

fairness -- resources should be used by all users fairly (perhaps in accordance with appropriate priority rules):

deadlock prevention == resources must be allocated so as to avoid deadlocks.

process between the pair of IMPs exchanging a message. Finally, Flow control and routing are distributed elegatithms which involve all the IMPs in the network. Each IMP makes local destatons about Global functions. The process of routing messages from source to destination involves all the IMPs in the network, in order that the best path for the message be shosen and agreed upon. Such distributed computations are quite different from conventional algorithms. They are not initialized, nor do they run to completion and hait. In a real sense, the ARPA routing calculation, flow gentrol techniques, and so on, have been in progress ever since 1969, because some part of the Network STHATE been running since then. These algorithms 479 gontinuously active progresses on a large number of different Processors. In fact, the number of processors and the interconnection between them is subject to change at any moment, They must run completely without human intervention. perform gontention resolution among bidders for shared resources, and they must do so without reces or deedlocks. (We have also dome to believe that it is essential to have a reset meshanism to "impossible" deadlocks and other conditions that may unteck result from hardware or saftware failures.)

2.1.2.2 Network Performance Goals

Packet-switching communications systems have two fundamental goals in the processing of data -- low delay and high throughput, Each message should be handled with a minimum of waiting time, and the total flow of data should be as large as possible. The difference between low delay and high throughput is important, what the network user wants is the completion of his data transmission in the shortest possible time. The time between transmission of the first bit and delivery of the first bit is a function of network delay, while the time between delivery of the first bit and delivery of the shortest bit is a function of network throughput. For interactive users with short messages, low delay is more important.

There is a fundamental tradeoff between low delay and high throughput, as in readily apparent in considering some of the mechanisms used to accomplish each goal. For low delay, a small packet size is necessary to cut transmission time, to improve the pipelining characteristics, and to shorten queueing latency at each mode; furthermore, short queues are desirable. For high throughput, a large packet size is necessary to degreese the circuit overhead in bits per second and the processing overhead per bit. That is, long packets increase the effective circuit

bendwidth and nodel processing bandwidth. Also, long queues may be necessary to provide sufficient buffering for full circuit utilization. Therefore, the network may need to employ separate mechanisms if it is to provide low delay for some users and high throughput for others.

To these two goals one must add two other equally important goals, which apply to message processing and to the operation of the network as a whole, First, the network should be cost-effective. Individual message service should have a reasonable dost as measured in terms of utilization of network resources; further, the network facilities, primarily the node computers and the circuits, should be utilized in a cost-effective way. Secondly, the network should be reliable, Messages accepted by the network should be delivered to the destination with a high probability of success. And the network as a whole should be a robust computer dommunications service; fault-etolerant, and able to function in the face of node of circuit feilures.

In summary, delay, throughput, reliability, and cost are the four criteria upon which packetwawitching network designs should be evaluated and compared. Further, it is the combined performance in all four areas which counts. For instance, poor

delay and throughput characteristics may be too big a price to pay for "perfect" reliability. What are the primary issues regarding the cost of computer networks? There are two kinds of costs considered heres the cost of the actual network components, and the cost of the use of the network. The first cost is a measure of the expense of connecting soms component to the network, and the second is a measure of the expense of utilizing the resources of that component. We are conserned here with outlining the effects of these costs on the balance among the various system paremeters of the network.

Law. East for Natural Compactivity. A major gonalderation in the cost of the ARPANET is the cost of connectivity. The basis variebles are sircuit costs and node costs; the values for the ARPANET are shown in Tables 2-1 and 2-2.

Bendwidth	Termination Cost	Line Cost 3/Month/Mile	
Kb#/806	3/Month		
4.6	650	. 40	
19.2	850 2,50		
50	850	5.00	

230.4

1350

10,00

Table 2m1 ARPANET Line Costs

Machine Type	Cost	Configuration (Kbs Tetal)
516 IMP	\$100K	up to 4 hosts (400)
		up to 5 lines (888)
		up to 7 devices tatel
316 IMP	\$50K	up to 4 hosts (300)
		up to 5 lines (600)
		up to 7 devices total
316 TIP	\$100K	up to 2 hosts (388)
		up to 3 lines (488)
		up to 64 terminale (100)

Table 2-2 ARPANET Node Coats

Several implications follow from those figures:

It The levout problem for networks is an important one, since reductions in the total number of circuit miles in a network represents substantial doller savings.

- 2. Low cost domes at the expense of both low delay and high throughput if lower line speeds are chosen.
- 3. Low dest also domes at the expense of low delay and high throughput if fewer lines are used, since more traffic is forced to use each line.
- 4. As a consequence of these points, some network designers, notably the Network Analysis Corporation, attempt to find a reletively low dost network layout and then test it by simulation to determine if average massage delay is below some threshold and the throughput obtained (when ell nodes send to all other nodes) is above some threshold.

Law Cost Top Nature & Use. A second papert of the cost of networks is the cost of the use of the network, which may be counted in one of several ways:

- 1. 3/bit, \$/character, \$/packet, or \$/message costs for the shipping of data through the network.
- 2. These sharges may be reted per mile or per hop or may be distance-independent.
- 3. There may be different grades of service at different costs.

These gosts are the translation into dollars of the utilization rates for various network resources. These resources can be detailogued as follows:

- 1. Time bandwidth
- 2. node processor bendwidth
- 1. node storege

Next we present some of the primary issues regarding the reliability of computer networks. In parallel with the discussion of network cost above, there are two basis topics examined here: the reliability of the network connections themselves, and the reliability of the network services. In the first instance, we are interested in how reliable the network components are. In the second case, the subject is the reliability of the use of the network facilities for data transmission. Admin, we are focusing in this section on the broad design issues which effect network performance as a whole. In section 2.3 on routing, we examine closely the effects that routing algorithms can have an network reliability.

Himi. Reliability. et. Metwerk. Conceptivity. One ean attempt to minimize the probability that a line or mode will be inoperative, and thus to reduce the probability that a mode connot dommunicate with the rest of the network. In practice,

this is a matter of maintaining a high magn time between failures and a low mean time to repair. One can also measure reliability in terms of the number of nodes and/or lines necessary to disconnect the network, taking into consideration the probability of the verious events. Alternatively, one can consider the size of the components of the disconnected network, or the fraction of node pairs not connected by any network path.

A different level of splution to the problem of network reliability is redundant network design, primarily in the layout of the circuits connecting the nodes. In this way, a network gan be constructed which is much more reliable as a whole than any one component. In the ARPANET, a design constraint has been that an IMP must be connected to the network by at least two circuits, so that the probability that an IMP cannot communicate with the network is very small. Of course, this principle pan be applied to other components in the network as well, The node computers can be backed up with siternate computers, or may be equipped with redundant interfaces and processors. The host computers size may wish to be connected to the network at more than one point by means of separate communications facilities.

High Reliability at Matwark Use. The reliability of the message programing on be message in terms of the persentage of

messages delivered, or the detected error rate, or the undetected error rate. Measures to improve the reliability of measure processing range from error detecting and correcting hardware to redundancy in software to backup message storage. The casts of these approaches are multiple; they add to the complexity of the and may degrade its performance, in addition to svstem. representing a dollars goat. In general, the communications subnetwork becomes more costly as these measures of reliability are improved. At some point, it becomes appropriate to pass the cost of these improvements on to the user. However, a minimum level of reliability is negestary for the operation of the communications subnetwork. Guarantees concerning a grade of mervice better then this minimum level might reasonably gost MOPS.

The Iradecti hetween Law Cent and Wigh Reliability. It is clear that there is a natural tradectif between lowecost networks and highereliability networks. This tradectif exists in building either spans networks or highly-gennected networks, and in providing special mechanisms to ensure the reliable transmission of data or shousing not to implement such sefectuards. In short, the price for reliability must be paid somewhere, either in the actual cost of constructing and maintaining the network, or in

the cost to the user of unreliable network service. It is difficult to generalize, but it may often be the dess that a lowedest network without sufficient measures for reliability may prove more costly to use in the long run than a network with a higher dost for higher network reliability. Stated differently, it may be cheaper to build a network to be fault-tolerent, redundent, and errer-detecting than to build these measures into each user process that communicates with the network.

2,1,2,3 Intended Functions

One of the most important factors in determining the overall design of any computer communications system is the set of functions and applications which must be supported by the system. In the case of the ARPA computer network, many different factors applied to make the design of a pioneering new approach to computer gommunications a necessary step. Foremost among these requirements was the need to connect terminals to computers with a very low response time. This was necessary in order to permit remote terminals to use time-sharing demputers as if they were connected directly to them. Other aspects computer environment are also worth noting. At the time the ARPANET was first built in 1969, there were no stendard computer communications protogols of any kind. Furthermore, connections between computer mainframes were extremely rare and usually special purpose. Very few interconnections of computers had ever been ettempted, prestigally none between computers of different manufacturers.

Thus, the ARPANET had to be greated with brand new technical approaches. The congept of a subnetwork of packetsswitching nodes, separate from the host computers, was a new idea. These computers were dedicated to the Job of communications, froming

the host computers for applications processing. Furthermore, these node computers needed to be very fast so that packets could be routed through the network through many intermediate peckets witching nodes in less than a second. Finally, it was shvistoned that the nodes should be sole to support relatively high amounts of traffic since the intended uses of the network included ramote job entry and file transfers, as well as interestive traffic. These high bendwidth applications were necessary in order to meet ARPAFs goal of resource sharing among the various ARPA-sponsored research centers.

Another important function to be performed by the communications subnetwork was interfecting between computers and terminals of widely varying types. Thus, loobit minisomputers should be able to communicate with 32-bit and 34-bit important computers and every computer should be able to communicate with every kind of terminal in use at the various sites. This formidable problem necessitated the design of a set of standard interfaces from host computers and terminals to the network and a set of higherential communications protocols for semmunication among the various types of equipment.

The intended functions of the ARPANET can be summarised as follows: general purpose intercommunication between a wide

variety of host computers and terminals with a minimum impact on those equipments and users and providing for a high level of performance.

2,1,3 Fundamental Design Choicee

We believe there are gix major areas in which the key choices must be made in designing a packet-switching networks

- extransmission Facilities
- -- Switching Technology
- -- Topological Structure
- **Communications Equipment
- Pecomunications Protogola
- seinterface Characteristics

2.1.3.1 Transmission Facilities

We next consider some of the important characteristics of the circuits used in the ARPANET.

The bandwidth of the network circuits is their most important characteristic. It defines the trafficecerrying

depectly of the network, both in the addregate and between any given source and destination. What is less obvious is that the bendwidth (and hence the time to glock a packet out onto the line) may be the main feator determining the trensit delays in the network. The minimum delay through the network depends mainly on direcult rates and lengths, and additional delays are lengtly accounted for by queueing delay, which is directly proportional to sincult bendwidth. These two factors lead to the general observation that the feater the network lines, the longer the packet can be, since long peakets have less overhead and permit higher throughput, while the added delay due to length is less important at high directly rates. In addition, more packet and message buffering is required when higher speed circuits are used.

The major effect of circuits with appreciable delay in a system requiring packet acknowledgment is that they require more buffering in the nodes to keep them fully loaded. That is, the node must meintain more packets in flight at once over a circuit with longer delay. This effect may be so lange (a circuit using a satellite has a delay of a quarter of a second) as to require significantly more memory in the nodes. This memory is needed at the nodes connected to the circuit to permit sufficient packet

buffering for nodesta-node transmission using the circuit. The subtle point is that additional buffering is also required at all nodes in the network that may need to maintain high source-to-destination rates over network paths which include this circuit. If they are to provide maximum throughput, they need sufficient message buffering to keep the entire network path fully loaded.

Traditionally, the telephone carriers have quoted error rates in the following menner; "No more than an average of 1 bit in 10 to the 6th bits in error; This definition is not entirely adequate for packet switching, though it may be for continuous transmission. For packet switching, the average bit error nate is less interesting than the average packet error rate (packets with one or more bits in error). For example, ten bits in error in every tenth packet is a 18% packet error rate, while one bit in error in every packet is a 188% packet error rate, yet the two desse have the same bit error rate.

An example of an ecceptable statement of error performance would be as follows:

The direuit operates in two modes, Mode is no sontinuous sequence of packet errors longer than

two seconds, with the average packet error rate

then one in a thousand. Mode 2: a continuous sequence of errors longer than two sequence with the following frequency distribution:

- > 2 seconds no more often than once per day
- > 1 minute no more often than once per week
- > 15 minutes no more often than once per month
- > 1 hour no more often than once per 3 months
- > 6 hours no more often then once per year
- > 1 day

While the figures above may seem too stringent in practice, the mode 1 bit error rate is actually quite lex sempered to conventional standards. In any case, these are the kinds of behavior descriptions needed for intelligent design of packet-switching network error control procedures. Therefore, it is important that the carriers begin to provide such descriptions.

The packet error rate of a circuit has two main effects, first, if the rate is high enough, it can degrade the effective circuit bandwidth by forcing the retransmission of many packets.

While this is basically a problem for the darrier to repair, the network nodes must recognize this gondition and decide whether or not to continue to use the direuit. This is a tradeoff between reduced throughput with the direuit and increased delay and less natwork connectivity without it. Before the direuit can be used, it must be working in both directions for peakets and for control information like routing and adknowledgments, and with a sufficiently low peaket error rate.

The second effect of the error rate is present even for relatively low error rates. It is necessary to build a very good erroredatection system so that the users of the network do not see errors more often then some specified extremely low frequency. That is, the network should detest enough errors so that the effective network error rate is at least an order of magnitude less than the host error or failure rate. A usual technique here is a cyclic redundancy check on each packet. This sheeksum should be chosen carefully; to first order, its size does not depend on packet length and it should be quite large, for example 24 bits for \$800000 lines and 32 bits for multi-megabit lines or lines with high error rates. (Note: essuming that the probability of packet error is proportional to the product of packet length and bit error rate, the checksum

length should be proportional to the log of the product of the desired time between undetected errors, the bit error rate, and the total bendwidth of all network circuits.)

2.1.3.2 Switching Technology

designing the ARPANET one of the key questions was how the subscriber equipment (hosts and terminals) would share the network transmission feetlities. Verious possibilities exists There are two general kinds of switching methods: point=4o=point mathods, such as multiplexing and switching, and multipoint or broadcast methods, such as polling and contention. In the first class, multiplexing, there is frequency division multiplexing, time division muitiplexing, and statistical time division multiplexing. At the time the ARPANET was designed, no off-the-shelf technology for multiplexing such diverse sets of subscribers into a common network was available. The key problem incompatibility of electrical and characteristics of the verious devices, Since multiplexing is inherently a transparent technique, no conversion would have been possible. With various switching technologies such as circuit switching, message switching, or packet switching, code conversion, data conversion, and spead conversion are possible, Circuit switching involves a dedigated physical path between the

communicating aubscribers, such as in the telephone system. This would have been inappropriate for the ARPANET due to the 1000 delays. Message switching usually involves a storewend forward process in which each node forwards the message to the next node which stores the message on the disk for later forward transmission. Such systems often have delays of many minutes. Packet switching was invented for use on the ARPANKT and can be seen as a logical outgrowth of message switching since it does not use secondary storage. The packets of each message are routed independently on a storesandsforward basis which allows pipelining of traffig through the network giving high efficienties while permitting low delays. Multipoint techniques ingluding polling were not used on the ARPANET primarily because perpenderance of asvachronous Charactermates time al the the verious ARPANET sites. terminals at A180, most of the ARPANET sites did not have computers which appeared their terminals on a polled basis. In fact, the absence of synchronous 1 newstweetime polled terminals from the ARPANET environment is quite an important aspect of the overall ARPANET design;

2.1.3.3 Topological Structure

The subject of network topology is a complex one, and we limit ourselves here to a few general observations. A

distributed topology was ghosen for the ARPANET as the legat=cost network design which would meet delay and throughput requirements.

problem of designing a general distributed network represented an important research problem in its own right ARPA contracted with Network Amelysis Corporation (NAC) to develop tools for designing such networks. There are many subproblems which can be identified in topological design, including node location, link capacity assignment, traffic flow assignment to links, combined sapesity and flow assignment, and finelly, the most general problem, topology sepecity and flow assignment. At the time that the ARPANET was designed, the only algorithms in existence tepelogical WOFE for **bestralized** Thus, NAC had to develop algorithms which found the minimum-cost network layout satisfying a cortain lavel of traffic depactty and minimizing delay given substriber locations, traffic flows, network element costs, and possible switch locations, This they did by a veriety of different kinds of algorithms, which cannot be discussed in detail here. A femily of procedures has been developed over the years which have been used to design optimal and near optimal connections for the ARPANET subscribers and which have been used to measure how close the actual

comes to cost effectiveness. The actual ARPANET configuration cannot be restructured every time a new subscriber needs to be connected. Therefore, it is always somewhat less then optimals. The studies by NAC have shown ARPANET topologies to be within a few percent of the best gost which could be expected.

The connectivity of the ARPANET modes was chosen to be relatively uniform, usually two or three. It is obvious that nodes with only a single line are to be evolded for reliability considerations. Nodes with many girouits also present reliability problem since they remove so much network connectivity when they are down. We also feel that the direction for future evolution of network geometries will be towards a "dentrol office" kind of layout with relatively fewer modes and with a high famein of nearby hosts and terminals. This tendency will become more pronounced as higher reliability in the node somputer begomes possible, even for large systems. One reason that we favor this approach is that a large node computer presents an ingressed opportunity for shared use of the resources (progessor and memory) among meny different devices a much more efficient and leading. to costueffective implementation. This trend will meen that in the future, even more than now, a key goet of network topology will be the

ultimate data connection to the user (host or terminal), who may be far from the dentral office. Concentrators and multiplexors have been the traditional solution; in packateswitching networks, a small node somputer should fill this function. In sonsitusion, we see flexibility and extensibility as two key requirements for the node computer. These fectors together with increasing performance and famely requirements imply a very high reliability extended as well.

2.1.3.4 Communications Equipment

The design of a large-scale computer communications network such as the ARPANET requires a set of communications equipment of several different types. In the case of the ARPANET which was designed at a time when the communications equipment industry was in its infancy, there were very few offethershelf communications processors to choose from. Today, the list is much londers

-- front-end processors

--multiplexors

-- line concentrators

-wterminal sontrollers (remote concentrators)

--eireuit awitches

-- message switches

Pepacket switches

senetwork manitoring centers

decision was made to construct the ARPANET using packet=Rwitching nodes. These nodes were designed from standard minicomputer base with additions for high speed interfaces to the wide band circuits and to host computers, Leter, the network grew to include terminal controllers which were merged with these packet switches. At an early stage in the network's development, the need for a network control center was clear and the development of the NCC hose was also initiated. The decision was made in most cases to locate the IMPs and TIPs near to the ARPA-exponsored fasilities. This eliminated the need in most cases for local access networks of any great complexity: Therefore, multiplexore and concentrators are not used in the ARPANET; all terminals and hosts connect directly to switching nodes. One exception to this rule has been the sonnection of some host computers over relatively long distances using the very distant host (VDH) interface which is simply a provision for wide band of roult connection between host computers and the ARPANET,

noteworthy development over the past several years has been the proliferation of frontwend processors in order to connect host computers to the ARPANET. This step was not the original intention of the ARPANET designers who felt that host computers should connect directly to the ARPANET to take full advantage of the power and flexibility of the network. subsequent experience has shown that the complexity of connecting some host computers directly to the ARPANET resulted in very expensive programming lobs with reduced efficiency within the At the same time as the trend towards frontwend processors has become evident in the ARPANET, these devices have also adhieved widewspread acceptance in the commercial market place. It is during this period of time that ISM introduced the \$784 and \$785 communications processors. (Of course, as we have noted earlier, the influence of the ARPANET has been seen in the introduction in the commercial market of packetsswitching nodes, network manitoring contert, and network-oriented controllers).

2.1.3.5 Communications Protocols

Communications protocols are designed to provide for the orderly exchange of data between computing equipment. They make possible logical communications over a given physical paths. Protocols operate in three basic modes:

- 1. establishing nedessary conventions, such as speeds and formats
- 2. establishing standard communications paths, including addressing, sequencing, error control, etc.
- 3. establishing the standard data element to exphange, sugh

The ARPANET protogols fit into this definition and have many other characteristics in sommon with other communications protocols. For instance, they are built on the principle of handwshaking, in which the sender and receiver exchange messages in pairs to insure the correct receipt of each message. The ARPANET protocols were also developed using the concept of a standard so that n senders and m receivers each convert to a standard protocol which requires n + m implementations instead of adopting pair-wise ad her communications conventions which would require n X m protocol implementations. In the ARPANET of protogol has been defined for communication between every pair of similar progesses. Likewise, an interfere has been defined for communication between each adjacent pair of dissimilar processes. Thus, there is a opotogo! between edlegent IMPs, between the source IMP and the destination IMP, between the source host and the destination host, between the sending terminal and the destination host, and so on, Also, there are interfaces between the IMP and the host, between the host operating system and the user program, ats, As is dommon predice in many computer systems, these different protogols have been designed modularly and in a hierarchical structure. This makes possible an independent design effort at the several levels and facilitates growth and ghonge within each level on an independent basis; Finally, it simplifies the specification and implementation of the higher level protogols since they can use the functions of the lawer level protogols instead of remimplementing them at the higher level.

2,1,3,6 Interface Characteristics

Each IMP will service up to four hosts whose cable distances, a from the IMP are less than 2000 feat. For greater distances, a modem shannel must be used. This latter type of host connection is termed a Very Distant Host (VDH).

Connecting on IMP to a wide veriety of different local hosts, however, requires, a hardwere interface, some part of which must be sustem tellored to each Host. It was decided, therefore, to partition the interface such that a standard

portion would be built into the IMP, and would be identical for all Hosts, while a special portion of the interface would be unique to each Host. The interface is designed to allow messages to flow in both directions at once. A bitmaerial interface was designed partly because it required fewer lines for electrical interfacing and was, therefore, less expensive, and partly to apprompt described appromption the variety of word lengths in the different Host computers. The bit rate requirement on the Host line is sufficiently lew that perallel transfers are not necessary.

The host interface operates asynghronously, each date bit being passed across the interface via a Ready for Next Bit/There's Your Bit handshake procedure. This technique permits the bit rate to adjust to the rate of the slower member of the sain and allows necessary interruptions, when words must be stored into or retrieved from memory. The IMP introduces a presedjusted delay between bits that limits the meximum data rate; at present, this delay is set to 10 migrosecends. Any delay introduced by the Hest in the handshake procedure further slows the rate below this 100 Kbe meximum.

The bandwidth and reliability of the heat gennection to the ARPANET are key parameters. The lesues in shooting the bandwidth

of the host connections are similar to those for the network In addition to establishing an upper bound on the hosts throughput, the rate is also an important factor in delay. The delay to send or receive a long message over a relatively alow host connection may be comparable in magnitude to the network round trip time. To eliminate this problem, and elso to allow high peak throughput rates, the host connection bandwidth should be as high possible (within the 8.0 cost-offectiveness), even higher than the average host throughput would indicate. By the same argument given above for pecket size, a higher speed host connection ellows the use of a lenger message with issu overhead and host processing per bit and therefore greater efficiency.

Seption 2

The reliability of the host connection is an important espect of the network design; several points are worth noting; First, the connection should have a packet error rate which is at least as low as the network circuits. This can be accomplished by a highly reliable direct connection locally or by error-detection and retransmission. The use of error control procedures implies that the host-node transmission procedures recemble the node-node transmission procedures which are discussed in a later section. Second, if the host application

requires extremely high reliability, a host-tonhoot date shecksum and message sequence check are both useful for detecting infrequent network feltures. Third, if the host requires uninterrupted network service, and the host is reliable enough itself to justify such service, multiple connections of the host to various nodes can improve the availability of the network, This option complisates metters for the source-to-destination transmission procedures in the nodes (e.g., sequencine) single there may be more than one possible destination node serving the host.

2.2 Topological Design

We next address the topic of topological design of the ARPANET.

2.2.1 Introduction

design of the topological levout of a computer communications system is an important design problem. There are many different possibilities for the topological atrusture of a network ranging from the aimple star configuration through various forms of hisrarchical networks to fully connected systems and distributed tepologies. The topological design problem is to find the minimum cost network levout given verious requirements such as subseriber locations, traffic flows, and performance These constraints are as variable as the computer sonstreints. communications systems involved. Typical performence constraints include average delay, peak throughput, and average reliability; These performance constraints affect a number of the choices of the topological design including the strautt speed, block size, number of links in the connections, equipment choices, and various protocol design choices es well.

2,2,2 Centralized Network Design

We begin with the problem of gentralized network design:
Here one is trying to obtain the least cost connection of a 491

of terminals to a center. The constraints require that particular directly have particular religibilities or throughputs and that an overall delay requirement be met. Without such constraints this problem reduces to the minimum spenning tree problem. There are two kinds of centralized natworks, the so-ceiled star natwork in which all terminals connect directly to the center, and the multipoint network in which particular communications singuits may be shared by more than one terminal. Two-level networks are also used in centralized designs; here the terminals connect to seme form of concentrators.

One-level centralized network design can be solved in two different ways. The optimal approach usually involves some kind of branch and bound technique in which the set of feesible solutions is partitioned into suggestively smaller subsets until a lower sost bound is achieved. Such techniques are diten too expensive to use for large scale networks, Neuristic approaches exist for the design of such networks. The EsqueWilliams algorithm proceeds by finding the node which is farthest from the denter, connecting it to the adjacent node which gives the greatest dost seving, and then repeating the process. Frimis algorithm selects the node closest to the center, connects the node closest to it, and repeats growing the gentralized natural node by node. Finally, the Kruskel algorithm works by picking the least gost link not yet used and inserting it into the

sentrelized network if it does not make a loop with some of the other links in the tree, It then repeats this process until a somplete tree is formed.

design of twowlevel centralized networks is more complicated since it is a hierarchical problem to determine how many concentrators to use, where they should be located, and which terminals should be gonnected to each concentrator. connection of the terminals to the congentrators may be achieved in a star or multipoint feebion as may be the connection of the soncentrators to the center. Here optimal solutions are not often employed and heuristiks wimilar to those for premievel networks are used. For instance, the Add algorithm, which is a form of "steepest descent", proceeds by connecting all the terminals to the center and then adding the concentrator which seven the most cost and repeating the process. The soughlied Drop algorithm is the inverse process. It connects all terminals to the least gost concentrator for that terminal and then removes the most costly concentrator, reconnecting the terminals to other soncentrators. The process is then repeated until a lower bound is achieved.

Centralized networks which meet performance constraints as well as cost constraints can be designed using these techniques. Furthermore, these algorithms can be used in the design of the

addess area network in large scale distributed systems as wells. To consider such systems, we begin by describing the design of the backbone net for a large decentralised network.

2.2.3 Backbane Network Design

The general problem of designing a distributed backbone for a computer communications network can be considered to be one or more of the following five subproblems in optimizations

- a. Node Location w Locate the switching nodes in the computer network.
- b. Capacity Assignment = Given the traffic flows and network topology, minimize average message delay with respect to channel depocity.
- g. Flow Assignment Given link departies and network topology, minimize everage delay with respect to traffia.
- d. Capacity and Flow Assignment Problem = Given the network topology, minimize the average delay with respect to channel capacities and traffic.
- e. The Topology, Capability and Flow Agaignment Problem = Minimize average delay with respect to network topology, capacity, and traific.

The capacity assignment problem can be solved exactly when channel costs are a linear function of bendwidth. Other solutions can be obtained for channel sosts which are logarithmic or power law costs. For stepsfunction costs, which are more common in practice, heuristic splutions are necessary.

The flow assignment problem sen be solved by the flow deviction method which is an optimal procedure resulting in an alternate routing flow. It begins by astablishing a feasible starting flow of traffic through the network. It then solves the shortest route flow problem by assigning lengths to each of the brenches associated with the increase in delay due to additional traffic on that link. The algorithm iterates by assigning more traffic to paths which represent the least additional delay paths to which some flow can be deviated.

When these two problems are combined to form the capacity flow assignment problem, globally eptimal solutions are no longer possible. Instead, least optimum solutions are obtained by beginning with a fessible starting flow calculating optimum capacity assignment under linearized costs, carrying out the flow deviation algorithm to find the optimum flows, repeating the sapecity assignment problem for those new flows and continuing to iterate between these two solutions until a least minimum is found.

The final problem is to determine the topology as well as the depectry and flow assignment. In this case an iterative algorithm known as someove branch elimination can be used. It works by selecting an initial topology and then carrying out the capacity flow sesignment algorithm. From a suboptimal solution this algorithm picks some specific capacities and gonducts a final flow optimization by an application of the flow deviation algorithm. Then it repeats these steps for a number of feasible starting flows. Finally, it repeats this entire election for a number of initial topologies.

2.2.4 General Network Design

The overell design of a general computer sommunications network incorporates elements of each of the procedures that we have described so far. One bagins by ideating backbone nodes in the network. Then the beekbone switching network is laid out using techniques such as the concave branch elimination algorithm and subalgorithms for determining depactities traffic and topology. A third step is to identify against area elusters of subscribers and to connect these subscribers to besidens nodes by means of one-level or two-level logal dentralized networks. Then one can repeat the design of the backbone and against area network with different starting essumptions or constraints until a satisfactory network design is equipped. The finel two steps are

ARPANET Completion Report
Chapter III Section 2 Design and Implementation

to identify the optimal carrier offerings with which to construct the communications network and the optimal communications equipment with which to build the network.

2.3 Communications Equipment

In this section of our report we disques the many different types of equipment which are contained in the ARPANET including the circuits themselves, the IMPs, TIPs, the NCC, and various other pieces of equipment such as the PLI, the RJE, etc.

2.3.1 The Network Circuite

IMP in the ARPANET can connect to up to four telephone lines. Normelly, these are 50 Kb per second elecuits. lines are provided by the Bell System as part of their widebend deta services. A 305-type data set is available as part of the terminal equipment for this facility. The physical connection to the local toll office is a twisted pair which resides in a cable containing about 1868 other twisted pairs. The 14ne Long haul elecuits are used to connect distant equalized. In that gase, incoming data at the logal toll office logations. is frequency-multiplexed and transmitted wideband over the long heut elreuite as an enelog signet with no coding. Approximately thirds of the long heal connections in the Bell System use microwave relay stations and the remainder use repeatered social sables. The 50 Kb per segand group feeflity can be operated in a synchronous or a nonsynchronous mode. The ARPANET uses the wideband service in the synchronous mode,

The 303-type wideband data station has been engineered seconding to the Bell System specifications to achieve a bit error probability on the order of one in 100,000. Errors occur for two basic reasons, impulse noise and momentary dropouts. Impulse neise seems to arise from strong field interferences such as those gaused by large gurrents which flow in switching senters. An even more important cause of arrors is momentary dropouts, These appear primarily on the long haul circuits which are known to be subject to gonditions which result in the occasional presence of a low signal-townoise ratio circuits. These conditions are generally referred to as "fades." The Bell System monitors their channels for fading conditions and provides for diversity in the event of a fade. Both types of errors tend to cause bit errors which occur in bursts.

Published experimental data on the nature of errors on leased private lines in the telephone plants indicates that an everage of one block of 511 bits in a few thousand blocks will not contain errors at all. Each block which contains errors will typically have errors occurring in bursts, and a block containing errors is also likely to have neighboring blocks containing errors. The following essumptions were made for the purpose of obtaining estimates on line errors in the design of the ARPANETI On a telephone line approximately one block is in error every three minutes on the average. This corresponds to approximately

one block of 511 bits in spror every 1888 blocks. To first order, the design of the ARPANET assumed that the error occurrences in time on the wideband festifities were approximately the same as those on the telephone lines. This means that an error every three minutes on the wideband service corresponds to one block of 1888 bits in error approximately every 18,888 blocks on the average at 58 Kbs. This size means that the distribution of error bursts over the wideband service differs from that on the telephone lines, whereas the median burst length on the telephone lines approximately 4 bits. The median burst length over the wideband service is larger by a factor between two end tens.

One of the important design objectives for the ARPANET was to establish essentially error-free communications so that any system difficulties sould be effectively and rapidly isolated. The error control acheme was designed to provide an average time between undetected errors in the communications subnetwork on the order of at least one or two years. The time between undetected errors is a random variable and therefore can have values which are lower than the average. An average time between errors of one or two years allows a mergin of safety which constitutes a reasonable design specification. The estimates of the network designers were that a loabit perity check would give an average time between undetected errors on gash ling in the ARPANET of

between helf a year and five years. With a 24-bit perity check, the estimate was that the time between undetected errors on gash line would be between 10 and 180 years. In a network with between 10 and 180 lines this gives an everage time between undetected error for the estimates of at least one year.

There ere two basis kinds of error detection mechanisms, syciic parity checks and lateral and longitudinal parity checks, also delied iterative checks. Cyclic parity checksums have been employed in the detection of errors on telephane lines, whereas iterative checksums have been used on magnetic tapes. The ARPANET design was based on the use of a BCH gods generated by a 24-bit shift register corresponding to the generator polynomial

The resear for selecting BCH godes was that the distant structure of these godes is better understood then most end they exhibit uniformly wide gode word separation. This means they have proved to be estimated shortes in detecting errors in many applications.

After the first four IMPs were delivered in 1969, a series of phone line tests were run on the initial ARPANET. The first test progrem continuously transmitted short packets on each phone line using looped telephone directs in order to perform two-ways tests. A 27-hour test indicated approximately one packet per

20,000 in error. The distribution of error bursts was concentrated in the range between one and seven packets. Subsequent tests were performed with much larger packets and over other lines for long periods of time. No errors were ever detected in the date by the phone line test program which had not elready been detected by the cyclic checksum hardware.

In the third quarter of 1978 when the ARPANET included it IMPs, the network was first tested with a 258,4 Kb per second diravit. Subsequently, a number of such high speed sircuits have been introduced in the network primarily for short distance limbs. These circuits are controlled by the IMP program without any special changes in software or hardware.

In the summer of 1973 the Norway TIP was connected by a 9.6 Kb per second satellite circuit to the rest of the ARPANET in the U_aS_a . The low speed circuit required some shanges to the IMP software to allow for slower propagation times for routing and to allow for more peckets in flight over the satellite channel.

The use of a satellite channel in a pecket-switching environment introduces some important new constraints to the method used to allegate satellite capacity. For example, in the SPADE system, channels are ellocated to voice conversations which are usually several minutes long. The minimum time required to allocate a channel pair is approximately one quarter of a second,

which is also the minimum time required to deallocate the channels. Allowing for additional delays due to ellocation conflicts, the total ellocationedsellocation time is on the order of one second, an insignificant fraction of the conversation time.

In contrast, the duration of a data packet in the ARPANET is sumparable to the duration of an SPADE dilocation request burst; thus if SPADE data channels were switched on a packet basis, each channel would be idle for at least the minimum 15 second allocation-deallocation time and utilized for only about 182 seconds essuming maximum length peckets. This is a utilization of only 4%.

An important additional consideration for packets witching traffic is the different delay requirements for different traffic classes. From these considerations it is given that if only single-packet messages are to be sent, it is best simply to send each directly with a destination address in an unused slot rather than sending a only request and release before and after it. On the other hand, for multi-packet traffic it becomes more efficient to first and requests for slots.

Several different approaches to the allocation of a satultite channel on a dynamic basis have been developed. These have been based on multi-access and broadcast modes of packet

delivery vie a destination address in each burst making use of the broadcast feature in varying degrees to accomplish dynamic eliocation of the channel.

The rendom edgess method also known as ALOHA is a pure sontention scheme. Other systems have been developed aumbining TDMA with reservation=based dynamic allocations. Soth of these systems have been tested in the ARPANET satellite IMPs using an INTELSAT=IV 56 Kb per second digital channel.

The remifications of the use of satellite channels for network trunks are considerable. Any satellite channel represents a path of higher delay and also potentially much higher bendwidth. If the channel is used in broadcast mode, then the delay and bendwidth characteristics of the link are more difficult to estimate presisely. In general, it is best to send short interestive treffic over terrestrial lines and long bujk data transfers over satellite paths.

The introduction of satellites into the system has a direct impact on nodel buffering. First of all, the IMPs sonnected to the satellite need much more peaket buffering than those sonnected to land lines since meny more packets can be in flight. Another important requirement and one which is less obvious is that all IMPs in the network need more memory for endetoeend message buffering when a satellite is added not just those IMPs

directly directed connected to it. The IMPs may also require special adknowledgment achemes to keep treak of the lenge number of packets in flight as well as changes in the host-to-host protocol and in other procedures.

In summary, the ARPANET was developed (nitially on the basis of 50 Kb per segond terrestrial sirguits. In subsequent years, it has inserporated lower speed and higher speed directs and satellite channels all without major shange to the original design.

2.3.2 The IMP Hardware and Softwere

The first packet-switch in the ARPANET series was based on the Honeywell Sib computer. The second packet-switch in the series was developed in 1971, based on the Honeywell Sib (the successor to the Sib). This second maghine was developed primarily to provide a less expensive version of the original switch. The choice of the Sib as successor to the Sib permitted this second version to be developed with minimal reprogramming, although hardware "specials" in the system had to be reengineered.

The experience of operating the ARPANET produced many insights into just what sharesteristics make a maghine more or isse suitable for use in a packateawitch. On the besis of this experience, a new machine called the Pluribus, was designed and built with the specific requirements of a packateawitching node in mind from the outset.

In Section 2,3,2,1, we discuss a number of packet-switching network considerations and their influence on the design of the packet-switching computer itself. In the remainder of Section 2,3,2, we describe the hardware and coftware attructures of the 516/316 packet-switch. The Pluribus-based switch is described in Section 2,3,5.

2,3,2,1 Network Considerations

The speed of the packet-switch's processor is an important determinant of the throughput rates possible in the network (the effect of processor speed on delay is of much less concern single processing delays are typically less than a millisequad). The store-and-forward processing bendwidth of the processor can be computed by counting instructions in the inner program loop. The source-to-destination processing bendwidth gam be calculated in a similar fashion. These rates should be high enough so that the entire bandwidth of the network lines can be used, i.e., so that the node is not a bottleneck. Experience indicates that the instruction cycle time is the main factor in this bandwidth deliculation; complex or specialised instruction sets would not be particularly valuable because simple instructions make up most inner loops-set less this is true in the systems we have built;

A different aspect of the packet-switches processor which sen also affect throughput is its responsiveness. Because circuite are synchronous devices, they require service with very tight time constraints. If the switch does not notice that input has completed on a given circuit, and does not prepare for a new input within a given time, the next input arriving on that circuit will be lost. Similarly, on output the switch must be responsive in order to keep the circuits fully loaded. This

DRAFT

requirement suggests that some form of interrupt system or high-speed polling device is necessary to keep response latency low, and that the overhead of an operating system and task scheduler and dispatcher may be prohibitive. Finally, we note that the amount of time required by the switch to process input and output is most critical in determining the minimum packet size, since it is with packets of this size that the highest pecket arrival and departure rates (and thus processing requirements) can be observed. Of source, data buffering in the device interfaces can partially alleviate these problems,

The speed of memory may be a major determinant of progessor speed, thus affecting the switch bandwidth. An equally important consideration is memory speed for I/O transfers, singe the switch's overell bendwidth results from a division of total memory bandwidth based on some processing time for a given amount of I/O time. First, there is the question of whether the I/O transfers ast in a cycle-stealing fashion to slow the processor or whether memory is effectively multi-ported to allew concurrent use. Second, there is the issue of sontention for memory among the various symphronous I/O devices. In a worst-case scenario, it is possible for all the I/O devices to request a memory transfer at the same instant, which keeps memory continuously busy for some time interval. A key design parameter is the ratio of this time to the available deta buffering time of the least

tolerent I/O device. This ratio should be less than one, and may therefore determine how much I/O gen be connected to the node.

The size of the memory, neturally, is enother key parameter. In our experience, with a terrestrial network the switch program and associated data atrustures take up the majority of the storage: this may change with the introduction of high-speed satellite sirsuits. The remainder of the switches memory is devoted to buffering of two kinds; packet buffering between adjacent nedes, and message buffering between source and destinction nodes. These requirements can be satsuisted quite almply in each case as the product of the maximum data rete to be supported times the benes trio time Cfor returning. In large metworks it may be necessary to rely acknowledoment). on sophisticated compression techniques to ensure that tables for the routing algorithm, the sourcesto-destination transmission procedures, and so on, do not require excessive storage.

The speed of the I/O system has been touched upon above in relation to processor and memory bendwidth. Other factors worth noting are the internal constraints imposed by the I/O system itself-wits delay and bandwidth. A different dimension, and one that we have found to be insdequately designed by most manufacturers, is the flexibility and extensibility of the I/O system. Host menufacturers supply only a limited range of I/O system.

options (some of which may be too plow or too expensive to use), Further, only a limited number of each type can be connected. A packet=switch requires high performance from the I/O system, both in the number of connections and in their data rates.

There are other factors to consider in evaluating or designing a computer for the packet-switch function apart from performance in terms of bandwidth and delay. As we mentioned, extensibility in I/O is very important and was comparatively rare until recently; it is more common to find memory systems which can be expanded. Processor systems which can be expanded are not at all common, and yet processor bandwidth may be the limiting factor in some switch configurations. Without a modular approach allowing processing, memory, and I/O growth, the cost of the switch dan be driven quite high by large step functions in component cost.

A final aspect of switch computer architecture is its reliability, particularly for large systems with many lines and hosts. A failure of such a system has a large impact on network partormance. Computer manufacturers tend to cut gorners in order to compete on price and delivery schedules. The panelty for this practice is usually paid in the soin of lowered reliability. These issues of performance, cost, and reliability become critically important in large networks serving thousands of hosts and terminals on a 24-hour-se-day basis.

2.3.2.2 The 516/316 IMP Hardware

The 516 and 316 IMPs can be gonsidered to be the same mechine. The 316 derivative is less expensive, emailer, effectively 30% slower, and less tightly engineered than the original ruggedised 516s. Architecturally, however, the mechines are very similar. They share the following characteristics: 1) simple processes with 16ebit word length; 512eword pages, single accumulator, approximately 1 migrosecond cycle time, one index register, indefinite indirect addressing, power-fail, and sutcerestarts 2) 16K words of memory which can be conveniently addressed; 3) a multiplexed 16-shannel direct memory access unit for Input/Gutputs 4) a 16-level priority interrupt system; and 5) a Teletype for maintenance and local debugging.

To the besit meghine are added a number of special hardware features which suit it to the IMP Job, specifically:

Full-duplex. Eveckronous interfaces provided far 470 IMP to the communication lines. the These interfaces are slocked by the modems at the bit level. They treated by the program at the pagket level (that is, they are direct access devices which send and receive blocks of words from storage, and use interrupts to notify the program of the completion of each packet transfer), The labs of denerating interspecket (idle) "syne" characters, providing packet framing sharesters, fetshing successive words of a packet from memory via a memory channel, serializing and describilizing all characters to and from the modem, formulating and checking a 24-bit systic redundancy ghock, detecting overflow and formet errors, and signaling the completion interrupt are all handled by the interface.

Full-duples, asymphronous interfaces connect the IMP to one or more host computers. Connecting on IMP to a wide variety of different hosts requires a hardware interface. fome part which must be gustom tallored to each host. A standard portion of the interface is built into the IMP, identical for all hosts (it is a direct memory channel device like the modem interface). While a special portion must be uniquely designed for each different host type. A bit serial interface is used partiv because it is less expensive and partiv to assummedate conveniently the variety of word lengths in the different hast computers. The interface operator asynchronously, each data passed agross the interfede via a Ready For Next beind Bit/There's Your Bit hendshake procedure. This technique permits the bit rate to adjust to that of either member of the pair and allows for pauses when words must be stored into or retrieved from memory.

A relative time clock provides an interrupt every 25,6 microseconds for perferming all timesdependent tasks such as timing out packets for retransmission. Programmed counters are run off this clock to perform still lower frequency periodic lobs.

A wetchdog timer is provided which is normally held off (i.e., restanted) by the program. The timer is restanted every time a correct response to a (periodic) "trouble report" is returned from the met indicating that most of the basic system is operating properly. If the timer ever runs out, a failure of some sort is indicated and a reloadwandwrestart program is equivated.

A programmegenerated interrupt, known as the task interrupt, permits programmegenerated tasks to be queued and handled by the same prioritywordered, interrupt-driven program structure as hardware (I/O) generated tasks.

A progremereadable hardwarm-hald register identifies the number and configuration of the machine.

It is illuminating to consider the impact of some detailed differences between the \$16 and \$16 machines. First, although the \$16 besie eyels is .96 microsecond and the \$16 dyels is 1.6 microseconds, the resi trafficehandling depathlities of the IMP

are not in this simple ratio singe other factors effect it. the 516, I/O channel pointers are kept in memory, while in the 316 they are in hardware. The result is that seek I/O word transferred in the 516 takes about 4 migroseconds while in the 316 it takes 3,2 microseconds. Thus for the same amount of I/O traffid, the 516 has less real time left for program execution, but it does instruction faster. In fact, for the IMP task, although one must dong der the memory access time used for I/O transfers, the processing bandwidth tends to predominate. relative importance of processing time to I/O time, however, will vary for a number of reasons, For instance, packet processing is less costly per bit for longer packets than short, whereas I/O transfer time per bit does not thenge. Second, in the 516 the watchdog timer is a hardware timer; in the 316 the timer is a software dounter run off the relative time glock interrupt routing. Neither timer can really verify that all features of the program are operating properly. The 316 timer is vulnerable to a break in a tiny piece of the glock gode but, by using the relative time alock, it eliminates the need for a separate hardware timer. We make up for this apparent deficiency by checking in other places in the sode to see if the clock code is running, and if it is not, a reload is initiated, Several other key data and control structures are tested in a similar fashion.

There are several other features which differ between the two machines but which have not turned out to be significant: 1) the 516 has a \$12-ward block of pretented memory, originally built in for protecting recovery and startup programs, while the 516 does not have this feature; 2) the helt instruction in the 516 can be inhibited (turned into a NOP) so that interrupts will always be able to force control to interrupt gerving routines while the 316 simply halts in all eases; and 3) the 516 has the option of halting or interrupting on power fail while the 316 is only able to interrupt.

2,3,2,3 The 516/316 IMP Beftware

Implementation of the IMPs required the development of a sophisticated operational computer program and the development of several auxiliary programs for hardware tests, program construction, and debugging. This section disquases the design of the operational program and describes the sayiliary software.

2.3.2.3.1 General Descriptions

As previously mentioned, the principal function of the IMP operational program is the processing of packets. This processing includes segmentation of Host messages into packets for routing and transmission, building of headers, receiving, routing and transmission of unacknowledged packets, reassembling of received packets into messages for transmission to the Host; and generating of RFNMs and adknowledgements. The program also monitors network status, gethers statistics, and performs oneline testing.

The entire program is composed of fifteen functionally distinct routines; each piece occupies no more than two or three pages of core (512 words per page). These routines communicate primarily through dommon registers residing in page zero of the machine which are directly addressable from all pages of memory. A typical map of core storage is shown in Figure 2-2. By use of

Figure 2-2: Typidal Map of Core Storage

a macro, code is "deniared" on tech phytical page at assembly time such that there is an integral number of buffers between the last word of code on one page and the first word of gode on the next. This technique eliminates all breakage (unusable except for part of one buffer on the very last page of core. Seven of the fifteen programs are directly involved in the flow of packets through the IMPs the task program performs the major portion of the packet processing, including the reassembly of Host messages; the modem programs (IMP=to=Modem and Modem=to=IMP) handle interrupts and the resetting of buffers for the Modem channels: the Host programs (IMP-to-Host and Host-to-IMP) handle interrupts and resetting of buffers for the Host channels. build packet headers during input, and construct allocation requests sent to the destination IMPs; the timeout program maintains a software clock, times out unused buffer allocations, reinitiates programs which have paused, and initiates routing somputations and other relatively infrequent events. A background loop contains the remaining major emetpota # (th and desis initialization, debugging, testing, statistics gathering, and Background programs also initiate RFNM allocation and tracing. other sequencing and control messages. After a brief description a) data structures, we will discuss packet processing in some detail.

DRAFT

Section 2 Design and Implementation

2,3,2,3,2 Date Structures

The major system data structures consist of buffers, queues and tables,

Buffer Stereds. The buffer storage space consists of about 52 fixed length buffers, each of which is used for storing a single packet. An unused buffer is chained onto a free buffer queue and is removed from this list when it is needed to store an incoming packet. A packet, once stored in a buffer, is never moved. After a packet has been successfully passed along to its Host or to another IMP, its buffer is returned to the free list; The buffer space is pertitioned in such a way that each process (store and forward traffic, Host traffic, etc.) is always guaranteed some buffers. For the sake of program speed and aimplicity, no attempt is made to retrieve the space wasted by partially filled buffers.

In handling store and forward traffig, all processing is on a parapacket basis. Further, elthough traffic to and from Hosts is composed of messages, the IMP converts to desling with packets; the Host transmits a message as a single unit but the IMP takes it one buffer at a time. As each buffer is filled, the program selects another buffer for input until the entire message has been provided for. These successive buffers will, in general, be scattered throughout the IMP's memory. An equivalent

inverse process occurs on output to the Host after ell packets of the message have arrived at the destination IHP. No attempt is made to dollect the peckets of a message into a dontiquous portion of IMP memory. A typical allocation of buffer apace in some storage is shown in Figure 2-2, as mentioned previously. IMPs with no Very Distant Host use the space on pages 35, 36 and 37 for buffer storage. All IMPs reserve the last 71 words for saving local date over reloads and for linkage to satellite and Terminal IMP programs.

The IMP program uses the following set of rules to ellocate the evaluable buffers to the various tasks requiring them:

- -- Each line must be able to get its shape of buffers for input. Double buffering is provided for input on each line, which permits all input traffic to be examined by the program. Thus, acknowledgements can always be processed, which frees buffers.
- buffers so that some lines may operate at full depactive.

 The number of buffers needed depends directly on line distance and line speed. The current limit is 1=1/2 times the number of channels per line, thus permitting 1=1/2 lines on the average to be operating at full depactty. Furthermore, each output line is guaranteed at

leest one buffer, thus permitting a low level of traffig on any line independent of congestion on other lines.

atoreda, including en overlep into the storemend-forward allocation dependent upon the number of lines. All IMP programma (except for modem input) must share this stored using a priority scheme to resolve contention and preglude "deadly embrace" = type stored lockups.

Buffers querently in use are either dedicated to an insoming or outgoing packet, chained on a quaye (or pointed to within a table) eveiting processing by the program, or being processed. Occasionally, a buffer may be simulteneously in several of these states, due to parallelism in the program. A "use count" is mainteined within each buffer, and only when this count goes to zero is the buffer put back on the free quaye.

Queues. There are three principal types of queues:

- -- Teski All routing pagkets, all paskets from the modems and all paskets received on Host channels, are placed on the tesk queue.
- -- Output) A separate output queue is constructed for each inter-IMP modem dirquit end each Host. Each modem output

queue is subdivided into a priority queue and a reguler message queue, which are serviced in that order. Eagh Most cutput queue is subdivided into a control message queue, a priority queue, and a regular message queue, which are also serviced in the indicated order.

== Reastembly; The reastembly queue contains those packets being reastembled into messages for the Host.

Tables. Tables in dore are identical for all IMPs, and their size is determined sither by fixed IMP peremeters of precessing capability and natwork performance considerations. In the former detegory, there are persIMP (67 word) tables for nouting and statistics data; perwhost (8 word) tables for queue pointer, status, and logal data required by rementrant Host code; and perwline (5 word) tables similar in function to the Host tables. In the latter setegory, all tables are pooled resources which are agguired and relinquished by IMP progresses depending on the activity required of them. These include transmit and receive message blocks, resessanbly blocks, trade blocks, transaction blocks, and initialization date.

The size of the initialization code and the essociated tables deserves mention. This was originally quite smalls However, as the network has grown and the IMP's capabilities have been expended, the amount of memory dedigated to initialization

has steadily grown. This is mainly due to the fact that the IMPs are no longer identically configured. An IMP may be required to handle a Very Distant Host, or TIP hardware, or five lines and two Hosts, or four Hosts and three lines, or a very high speed line, or a satellite link. As the physical permutations of the IMP have continued to inspease, the criterion followed has been that the program should be identical in all IMPs, allowing an IMP to reload its program from a neighboring IMP and providing other considerable adventages. However, maintaining only one version the program means that the program must rebuild itself during initialization to be the proper program to handle the particular physical configuration of the IMP. Furthermore, it must be able to turn itself back into its nominal form when it is reloaded meighbor. All of this takes tables and code. Unfortunately, the proliferation of IMP sonfigurations which has taken place was not foreseen; therefore, the program differences currently gannot be conveniently computed from a (mp) configuration key. Instead, the configuration irregularities must be explicitly tabled.

2,3,2,3,3 Packet Flow Through Major IMP Routines

Figure 2=3 is a schematic drawing of packet processing. The processing programs are described below. Facket flow may be followed by referring to Figure 2=3.

Figure 2-3: Pecket Flow and Processing

Haststevide routine (HeI) hendles messages being transmitted into the IMP from a local Host. The routine first accepts the leader to construct a header that is prefixed to each packet of the message. It then accepts the first packet and, if no allocation of space exists for the destination IMP, constructs a request for buffer eligeation, which it places on the task Single-packet messages are placed directly on the task quave regardless of allogation status and are held via a transaction block until either a RFNM or ellocation is returned. A returned RFNM releases the pagket. A returned allogation for the single-peaket message will cause petrenamission from the background loop. Requests for multipagket allogation are sent without actual message data. The request is recorded at the destination IMP and an allogation massage is returned via the background loop when space is available. A returned allocation seuses HeI to release the first packet with header to the tesk queue vie the progremmable task interrupt, Subsequent input is then accepted from the Host until end of message (EOM) occurs. The routing also tests a hardware trouble indicator and verifies the message format. The routine is serially reentrent and services all Hosts connected to the IMP.

The Modemetseims routine (MeI) handles inputs from the modems. This routine first sets up a new input buffer, normally obtained from the free list. If a buffer cannot be obtained, the

ARPANET Completion Report Chapter III

Section 2 Design and Implementation

received buffer is not acknowledged and is reused immediately. The discerded packet will be retransmitted by the distant IMP. The routine processes returning acknowledgements for previously transmitted packets and either releases the packets to the free list or signals their subsequent release to the IMP-to-Modem routine. The Mai routine then places the buffer on the end of the task queue and triggers the programmeble teak interrupt.

The TASK routine uses the header information to direct packets to their proper destination. The routine is driven by the task interrupt, which is set whenever a packet is put on the task quaue. The routine routes packets from the task quaue onto en output modem or Host queue determined from the routing alcorithm. If the peaket is for non-local delivery, the routine determines whether sufficient atore and forward buffer space is available. If not, buffers from modem lines are flushed and no subsequent eaknowledgement is returned by the IMP-sto-Modem routing. Normally, an adknowledgement is returned in the next Daloptuo packet over that modem 14mm. Packets from Hosts which gannot det store and forward apage are removed from the queue and replaced at a later time by the Hell routine,

If a pecket from a modem line is addressed for local delivery, its message number is sheaked to see whether a duplicate packet has been received. As mentioned previously,

sech IMP meintains for each connection a window of contiguous numbers which it will accept from the other side of the connection. Packets with outsoference numbers are considered duplicate and are discarded. The regains of a RFNM for the eldest message by the source Host permits the window to be moved up by one number.

Replies such as RFNMs or Dead Host messages are placed in transaction blocks. TASK then pokes the IMP=to=Host routine to initiate output to the Host.

Message packets for loss delivery are linked together with other packets of the same message number in a reassembly block, when a message is completely ressembled, the leading packet is linked to the appropriate Host output quaue for prosessing by the IMPetoeHost.

Indoming routing messages are processed to that outgoing routing messages and the routing directory immediately reflect any new information received. The task routine generates I-heardwyou messages in response as negestary to indicate to the neighbor receipt of the routing message.

IMPringMades (I-M). This routine transmits successive packets from the modem output queues and sends piggybacked acknowledgements for packets correctly received by the Modemato-IMP routine and accepted by the task routine.

IMPRINGER (I=H). This routine person messages to lode! Hosts and informs the background routine when a RFNM should be returned to the source Host.

Initialization and Baskszound Legg. The IMP program starts in an initialization section that builds the initial data structures, prepares for inputs from modem and Host channels, and resots all program switches to their nominal state. The program then falls into the background loop, which is an endlessly repeated series of lowerfority subroutines that are interrupted to handle normal traffic.

The programs in the IMP beakground loop perform a variety of functions: TTY is used to hendle the IMP Teletype traffigs DEBUG, to inspect or change IMP core memory: TRACE, to transmit collected information about trace packets; \$TATISTICS, to take and transmit network and IMP statistics; PARAMETER=CMANGE, to alter the values of selected IMP parameters; PACKET CORE, to transfer portions of core images via the network; and DISCARD, to throw away pagkets. Selected Hosts and IMPs, particularly the Network Control Center, will find it necessary or useful to communicate with one or more of these background loop programs. these progrems may send and receive messages from the So that network. they are treated as "fake Hosts.* Rather then duplicating portions of the large IMP-to-Host and Host-to-IMP routines, the background loop programs are treated as (f they were Hosts, and they can thereby utilize existing programs. The "For IMP" bit or the "From IMP" bit in the leader indicates that a given message is for or from a fake Host program in the IMP. Almost all of the background loop is devoted to running these programs.

The TTY program essembles characters from the Teletype into network messages and decodes network messages into characters for the Teletype. TTY's normal message destination is the DEBUG program at its own IMP; however, TTY can be made to communicate with any other IMP Teletype, any other IMP DEBUG program or any Host program with competible format.

The DEBUG program parmits the operational program to be inspected and changed. Although its normal message source is the TTY program at its ewn IMP, DEBUG will respond to a message of the correct format from any source. This program is normally inhibited from changing the operational IMP program; Network Control Center intervention is required to activate the program; full power.

The STATISTICS program goiledts measurements about network operation and periodically transmits them to a designated Hosts This program sands but does not receive messages. STATISTICS has a mechanism for collecting measurements over 10=second intervals

end for taking helfspecond anapshots of IMP queue lengths and routing tables. It can also generate artificial traffic to load the network.

The PACKET CORE program loads and dumps portions of its own IMP's dore memory, or acts as an intermediary in loading and dumping portions of the core memory belonging to a neighbor who is unable to communicate via the normal IMP-IMP protocol. The PACKET CORE facility allows for dissimilar machines to doexist as IMPs on the networks reloading and diagnostic dumping of a malfunctioning IMP can be done without the requirement that one of its neighbors be of the same machine type.

Other programs in the beckground loop drive loop; status lights and operate the parameter change routine. A thirty-two word perameter table controls the operation of the TRACE and STATISTICS programs and includes spares for expansion; the PARAMETER=CHANGE program accepts messages that change these peremeters.

Other routines, which send connection protocol messages, send incomplete transmission messages, send sliceations, return givebacks, send RFNMs, and retransmit single-packet messages also reside in the background programs. These routines are called Sack Hosts, However, these programs run in a slightly different menner than the fake Hosts in that they do not simulate the

Section 2 Design and Implementation

Heat/IMP channel hardware. They do not go through the Host/IMP gode at all, but rather put their messages directly on the teak queue. Nonetheless, the principle is the same.

Timent. The timeout routine is started every 25,6 maggeriled a featwhick timeout period) by a clock interrupt. The routine has two sections: the feat timeout routine which "wakes up" any Host or modem interrupt routine that has languished (for example, when the Host input routine gould not immediately start a new input because of a shortage in buffer spece); and the slow timeout routine which marks lines as alive or dead, updates the routing tebles, and does long term garbage collection of queues and other data structures. (For example, it protects the system from the cumulative effect of such failures as a lost packet of a multiple packet message, where buffers are tied up in message reassembly).

These two routines, Past and Slow, are executed so that fest timeout runs every glock tick (25,6 ms) and the slow timeout runs every 25th clock tick (648 ms). Although they run off a common interrupt, they are constructed to allow feet timeout to interrupt slow timeout should slow timeout not complete execution before the next timeout period. During garbage sollection, every table, most queues, and many states of the program are timed out. Thus, if an entry remains in a table abnormally long or if a

routine remains in a particular state for abnormally long, this entry or state is carbage-collected and the table or routine is returned to its initial or nominal state. In this way, abnormal conditions are not allowed to heng up the system indefinitely.

In addition to timing out verious states of the program, the timeout routine is used to awaken routines which have put themselves to sleep for a specified period. Typically these routines are weiting for some resource to become sveilable, and are written as de-routines with the timeout routine. When they are restarted by Timeout the test is made for the aveilability of the resource, followed by another delay if the resource is not yet evaluable.

2.3.2.3.4 Control Organization

It is characteristic of the IMP system that many of the main programs are entered both as subroutine dells from other programs and as interrupt dells from the hardware. The resulting control atrusture is shown in Figure 2=4. The programs are arranged in a priority order; dontrol passes upward in the shain whenever a hardware interrupt occurs or the surrent program decides that the time has dome to run a higher priority program, and control passes downward only when the higher priority programs are finished. No program may execute either itself or a lower priority program; however, a program may freely execute a higher

Figure 2=4: IMP Program Control Structure

priority programs. This rule is similar to the usual rules concerning priority interrupt routines.

In one important case, however, control must pass from a higher priority program a namely, from the several input routines to the task routine. For this special case, the computer hardware was modified to findlude a low-priority hardware interrupt that can be set by the program; when this interrupt has been honored (i.e., when all other interrupts have been serviced), the task routine is executed. Thus, control is directed where needed without vieleting the priority rules.

The practical implementation of priority control involves the setting of interrupt masks and emabling or inhibiting Masks are built during initialization. In general, routine is entered, either by hardwares when software-initiated interrupt, the entering mask registers and keys are seved. A mask for the new routine is set into the mask register, and the routine controls interrupts by executing INH or ENB commends. Therefore, HeI may inhibit interrupts by MwI for short periods of time during spitisal functions by using the INM. When the ENB command is executed, however, the mask bits for M=I will permit herdwere interrupts transferring control from HwI. Interrupt control is obviously extremely oritical and its use constitutes the most complex area of program operation,

DRAFT

Section 2 Design and Implementation

All interrupt levels (except modem to IMP) make use of interrupt entry and exit routines that save the interrupted atete variables on a stack. The six veriables saved ares the index register, the accumulator, the keys (overflow bit, memory mode, etg.), the checksum routine return address, the interrupt mask, and the interrupt return address. The checksum routine is therefore resentrant at each separate interrupt level and can be called with interrupts enabled. For a slight improvement in afficiency, the Modemeto-IMP routine does (to own saving and restoring of state variables.

Some routines must occasionally wait for long intervals of time, for example when the Host-to-IMP routine must wait for an ellocation from the destination IMP. Stopping the whole system would be intolerable. Therefore, should the need arise, such a routine is dismissed, and the timeout routine will later transfer control to the weiting routine.

The control structure and the partition of responsibility among various programs achieve the following timing goals:

- -- No program stops or delays the system while weiting for an event.
- -- The program gracefully adjusts to the situation where the machine becomes computerbound.

- The Modemato-IMP routine can deliver its current packet to the task routine before the next packet arrives and can always prepare for successive packet inputs on each line. This timing is critical because a slight delay here might require retransmission of the entire packet.
- be sent as fast as they can be accepted by the phone
- -- Necessary periodic progresses (in the timeout routine). are always permitted to run, and do not interfere with inputsoutput progresses.

2.3.3 The TIP Hardware and Software

The Termine) Interface Message Processor (TIP) provides a means for connecting up to 63 termine) devices to the ARPA Network. The terminal interface specification conforms to the EIA standard R\$232C, which permits direct connection to most data modems. In addition to full duplex, series data transmission, each of the 64 ports provides 4 programmestiable control lines and monitors 6 external status lines, these lines are useful in desling with modems or other compatible I/O devices. Data formatis Teletype compatible, that is, character oriented with start and stop bits. The TIP handles all routine operations of timing and sequencing. All line persmeters, such as speed and sharecter size, are program settable.

2.3.3.1 The TIP Herdware

The TIP is built eround a Moneywell He316 computer with 28K of core; It embodies a standard is-port multiplexed memory channel with priority interrupts and includes a Teletype for debugging and program reloading. Other feetures of the standard IMP also present are a registime glock, powereful and sutpersettent mechanisms, and a programmegenerated interrupt feeture. As in the standard IMP, interfaces are provided for connecting to highespeed (58-xilobit, 238,4-xilobit, etc.) modems as well as to Hosts.

Aside from the additional 12K of core memory, the primary hardware feature which distinguishes the TIP from a standard IMP is a Multimeline Controller (MLC) which allows for connection of terminals to the IMP. Any of the MLC lines may go to loggiterminals or via modems to remote terminals. As shown in Figure 2-5 the MLC gonsists of two portions, one a piece of central logic which handles the assembly and disassembly of characters and transfers them to and from memory, and the other a set of Line Interface Units (all identical except for a smell number of option lumpers) which synchronize reporting to individual data bits between the central logic and the terminal devices and provide for control status information to and from the modem of device. Line Interface Units may be physically incorporated one at a time as required.

Figure 2=5 shows the hardware configuration of the TIP. The upper part of the figure shows the essentially standard 516 IMP contained in a TIP which runs the standard IMP program. This program includes a test for the presence of a Multi-Line Controller. If an MLC is attached, the program estivates a terminal-handling program which controls the MLC.

The lower portion of Figure 2-5 shows the Multi-Line Controller. The MLC provides interfaces for up to sixty-four full duplex connections. These lines may be glocked either

Section 2 Design and Implementation

Figure 2-5 TIP Hardware Configuration

externally to the TIP, in which case they are called externally clocked, or internally by the TIP, and called asynchronous. Any of these connections or ports (as they will be referred to hereafter) may transmit and receive up to 19,2 Kbps when operating externally clocked. The bit rates for the internal clocks range from 75 band to 2400 band and are shown in Table

MLC CLOCK RATES

Table 243

2=3.

All of the directory specific to a single port is contained on a plugain module salled a Line Interface Unit or LIU. Although the input and the output for a given port are usually seafgned to the same devices they may so to completely separate

devices such as a card reader and a line printer. All such perameters as bit rate and character size are individually settable under program control for each port and are independent for the input and the output of a given port.

In addition to circuitry on the LIU, each port uses logic which is shared by all lines. This is called the common common logic contains the DMC interfeces, the The timing directory, and a 75-bit by 64-word memory. Each of the 64 device ports is assigned a unique memory location called the STATEWORD for that line. STATEWORD records the instantaneous state of data and control transmission on a line. Part of the memory is realized as a serial, pseudosdrym memory built from disculating shift registers. This 57=bit portion is shown schematically in figure 2.6. The remaining 18 bits of the STATEWORD are stored on a per-line besis on the individual LIUs. The pseudo-drum rotates once per 51,2 migroseconds which permits a complete revolution per bit at 19,2 kilobits per second. Thus the MLC handles one bit per revolution at 19.2 Kbps, one bit per 2 revolutions at 9.4 Kbps and so forth,

All external directory is logated on the Line Interface Unit, The LIU contains circuits which denorate signals which may not be time multiplexed, i.e., turned on and off as the paeudo-drum rotates. Such signals are the output data line and

Figure 2=6 Stateword Pseudo=drum Memory

modem control lines. The LIU contains a portion of the STATEWORD, including the output data synchronizers, the output data synchronizers, the output dock synchronizers, control signals for an I/O device or an optional modem, the input data synchronizer, and the input clock synchronizers.

The format for all data which will be used as input for a given port of the MLC must be either five, six, seven or eight-bit character. Each character must also be preceded by a start bit and followed by a stop bit. Essentially, characters must be compatible with standard Teletype format. The interfess voltage levels conform to EIA Specification 88232 and the interface connector is the EIA standard dataphone connector. In the standard configuration the TIP is sveliable with up to eixteen internally mounted modems of the 102, 201, or 202 variety.

For externally glocked date (that is, slock signal provided by the modem or the attached device), when the receive clock makes a negative transition, the input side of the device port will read the state of the device data line. For externally clocked output, the device port must change data in response to positive transitions of the external clock line. For data input where there is no external clock line to tell exactly when to sample the data ine, the eightephase bitesampling technique,

which employs one of the internally generated glock frequencies, is used. On internally glocked output, no genetraint is imposed on when a character may start; however, once a character has started and generated a start bit on the outgoing data line, it will proceed at the specified and welleddined rate. Onset of the next character may been no relationship to the present character. It is for this reason that this mode is normally sailed esynchronous.

For operation of the device control and status monitoring, ten lines are evaluable. There are four control lines which are set under program control; these conform to EIA Spec. R\$232, wherein a positive voltage is a functionatrue operation. There are also six status lines from the modem into the Multi-Line Controller. They may be read at any time under program control, and they are defined as for the control lines; positive voltage equals a functionatrue condition.

This control status information is communicated to and from the 316 domputer via the accumulator imputabutput bus, which has an average delay of 25 microseconds and a worstweese delay of 58 microseconds when used through the Multistine Controller. To these delays must be added any delays due to the operation of the TIP program. Consequently, this is not a highespeed data path and is used only for control information at a low rate.

The next section describes the manner in which serial bit streams from a device port are assembled into 5-bit sharesters and subsequently loaded into the 316 computer memory. The inverse operation of accepting 8-bit sharesters from the memory and disassembling them into bit streams for the verious ports is elso described.

The Multi-Line Controller communicates with the 316 via memory (i.e., DMC) channels, program-generated control signals, accumulator I/O instructions and priority interrupts. The program exchanges characters with the MLC through buffer tables in the memory, accessed via the DMC. One DMC channel is used for input and two for output. Two priority interrupt lines are used; one to indicate that the output buffer table has been emptied, and one to give a periodic interrupt which causes the program to check for input and for output requests (see below for details).

The program can set pertain bits in the MLC for each line. For every line's input section, the program sets cherecter size, line rate and whether or not the line is of high input priority, These bits are set through an instruction which is addressed to the controller and whose date word designates the input section and includes the line number as well as the specific parameter values. This is used to initialize each input line on startup.

The same instruction with the data word designating the output section is used to initialize the output lines. This command sets the output line speed to the appropriate value.

Instructions are provided for starting and atopping transfers by the controller on the DMC channels. For exemple, when the program wishes to look at the input table, it executes an Input Disable deusing the controller to inhibit further requests on the input channel. After providing a new buffer it executes an Input Enable which restarts the input transfers, Similar instructions are provided for controlling MLC access to a memory table where persion requests for output characters are reported to the program.

The Output table, which contains outgoing characters to terminals, is reestablished only when the table is completely empty (program-controlled overall stopping of output is not provided), so only the Output Enable is required.

Figure 2=7 shows the Input Data Flow, For each line, the MLC provides hardware storage for the character presently being shifted in and for the previous character. The serial bit stream enters a shift register called INDAT, when a character is fully exsembled, it is moved to a buffer, BINDAT, which is then marked as full. Each line has a private INDAT shifter and BINDAT buffer, During the period that a character is being shifted into

INDAT, the previous character in BINDAT must be unloaded into the Because of variations in character length and particularly in line appeds, the amount of urgency in emptying BINDAY varies from line to line. The input DMC channel is used to pump the characters from the BINDATs (nto the memory, If BINDATE were serviced solely on a first-comesfirst-served basis, an urgent high speed line gould get locked out by less urgent low speed lines. Worst gase memory appeasing is further gomplicated by higher priority DMC channels and by the fact that the 64 BINDATE are not randomly available but are seanned in sequence (one every 800 hs). Lines operating at under about 3000 bits/sec, even in the worst case with the logkout by all other lines and the expected interference from Hosts (essumed 100 Kb) and modemme (assumed 50 Kb), will get their BINDATs ampty in time for the next character. Faster lines can be locked out until too lete (i.e., until efter their BINDAT is needed for the ensuing character). To avoid this, two buffer registers are provided for entering characters into the single input DMC channel. lines use one of the buffers (FINBUF) and normal lines use the other (INBUF). Lines are marked in the MLC by the program as FAST or not.

Input characters are dymped into a tymble table. Each entry consists of the character, the line number (9=63), and a bit that indicates whether the input side of the line is in BREAK mode.

Section 2 Design and Implementation

Figure 2=7 Input Date Flow

When the progrem is ready to start accepting input, it establishes an input buffer (sets pointers) for the DMC input channel and executes an Input Enable, It then sets up character size, line rate, and "Fast" priority for each line in turn. The lines are now active; characters come into the controller and are stored by it into the input table in sore.

Periodically (every 3.3 milliseconds), a clock interrupt rouses a routine which checks to see whether or not any characters have been received. The buffer provided is sufficiently long that it should never be filled by the hardwere between periodic looks. In order to check for input, the program blocks input on that DMC channel and dauges the controller to rewrite the DMC pointers in the memory where the program den look at them. The program seves the pointer values and then resets the pointer end executes enother Input Enable, which reopens the door for further input.

Figure 208 shows the Output Deta Flow. The MLC contains a shift register and a character buffer for each output line. One of those, Outpat, contains the character presently being shifted onto the line; the other, BOUTDAT, contains the next character to go. As soon as the character in OUTDAT has been completely shifted onto the line, the character in BOUTDAT is moved into OUTDAT and in turn starts to be shifted out. A full character

time is thus evellable for fetching the next character from memory while ensuring that the line is kept busy. If a character is not evellable for output when needed, autput seases until a character becomes evellable. Since each character has a START and a STOP bit this does not interfere with character synchronization.

For output, the only line peremeter set by the program is line speed. Character size is directly determined by the program for each character put out. There are 10 bits evellable for holding an output character (in BOUTDAT and also in OUTDAT). The START bit is automatically erested by the hardware for each character sent to the line and thus the program can ignore this format requirement.

Output is hendled by the progrem through tebles sedessed by two DMC chennels. One of these, the OUTPUT teble, contains the outgoing cherecters. Since the herdware takes the cherecters from the table in sequence, care must be exercised to prevent trying to feed a cherecter to a port which cannot presently take it, thereby possibly blocking cherecters to more needly lines. (A cherecter can Jem the OUTPUT channel, sitting in OUTSUF until the appropriate lines SOUTDAT can accept it.) The program must thus cleverly feed characters to lines at the proper rates depending upon individual line speeds. To essist in this process

Seption 2 Design and Implementation

Figure 2=8 Output Data Flow

a second table, the GUTIN table, is used. Each time a line's BOUTDAT buffer empties, a request is entered into the OUTIN table which effectively says, "I den take another character now." (The entry in the table consists simply of the line number.) The program periodically checks the OUTIN table (handling it just like the INPUT table discussed above) and feeds characters into the OUTPUT table for needy lines. Since several lines may need to enter requests at about the same time and the DMG channel may not be able to suck up the requests as fast as the controller wents to generate them, a bit is provided for each port in the MLC to remember that a request needs to be entered. The bit is cleared when the port succeeds in entering its request.

A line will be kept busy if a new sharacter can go from the DUTPUT table into the BOUTDAT between the time BOUTDAT empties and the time the character has been shifted out to the line (nomplete with stop bit(s)). For slow speed lines there is adequate time for the request to be entered, the progrem to note the request and put a character into the OUTPUT table, and the character to get into BOUTDAT.

The program looks at the OUTIN table every 3.3 millimeands, Within that period of time only one character can have been transmitted on the lines whose character period (including START and STOP bits) is greater than 3.3 ms. For an 8-bit character

ARPANET Completion Report Chapter III

with one START and one STOP bit, this means a bit rate of about bits/second. Thus, lines of 3000 bits/second and under using characters with 8 data bits and 1 STOP bit gan only enter one request into the OUTIN table between program looks, and a new sharester will be put into the GUTPUT table, in response to this request, before the line "runk dry". Blockage of the output table could, it would appear, prevent the character from going aut but, if we ignore higher speed lines for the moment, any line with a request in the GUTIN table has an empty BGUTDAT and thus can accept the character "immediately." Actually, all lines sould dome due for characters at once and since DMC access requires something over 3.2 migroseconds and access to the SQUIDATS is sequentiel, it can take some time to get a character out to a given line. At worst, we can get out a character to some line about every 50 migroseconds and thus to det a character For any line requires at worst about 3 ms (64 x 50 migroseconds). For 8 bit/1 STOP bit characters, this means that the bit rate (8 3000 bits/second. Thus, lines operating at 3000 about bits/second or less can always get a character from the output table in time, even if all other lines are queued up ahead of it waiting for characters. As we have seen, the program will process requests at about that rate and thus lines of less than 3888 bits/second should operate with no degradations

Lines of greater than \$600 bits/second may not always be able to get out a character in time and thus for these lines, the characters must be planted sarefully into the output table by the program so that they come out frequently enough to keep the line going. By cleverly watching the OUTIN requests, the program can applicable the output character for high speed lines into the OUTPUT table so that they do not hold up these lines.

Note that all of these galquietions are dependent on sharecter length (including START and STOP bits) shrinks, the problem becomes more south as the gharacter rate increases for a constant line speed in bits/second.

2.3.3.2 The TIP Software

Because the terminals connected to a TIP gommunicate with Hosts at remote sites, the TIP, in addition to performing the IMP function, also acts as intermediary between the terminal and the distant Host. This means that network standards for format and protocol must be implemented in the TIP. One can thus think of the TIP software as containing both a very simple-minded mini-Host and a regular IMP program.

Figure 2=9 gives a simplified diagrammatic view of the program. The lower block marked "IMP" represents the usual IMP program. The two lines into and out of that block are logically

equivalent to input and output from a host. The gode conversion blocks are in fact surprisingly complex and include all of the material for dealing with diverse types of terminals.

Once connection to a remote Most is established, regular messages flow directly through the Input block and on through the IMP program. Returning responses dome in through the IMP, into the OUTPUT program where they are fed through the OUTPUT Code Conversion block to the terminal itself.

As the user types of the keyboard, characters go, via input code conversion, to the input block. Information for remote sites is formed into regular network messages and passes through the OR switch to the IMP program for transmittely. Command characters are fed off to the side to the command block where commands are decoded. The commands are then "performed" in thet they sither set some appropriate parameter or a flag which calls for later setion. An example of this is the LOGIN command. Such a command in fest triggers a complex network protocol procedure, the various steps of which are performed by the PROTOCOL block working in conjunction with the remote Host through the IMPs. As part of this process an appropriate special message will be sent to the terminal via the Special Messages block indicating the status (success, failure, stall of the procedure,

Figure 2=9 Block Diagram of TIP Program

ARPANET Completion Report Chapter III

DRAFT Section 2 Design and Implementation

2.3.3.3. THE TIP Command Format

The user at a terminal will at various times be talking directly to his TIP instead of to the remote Hosts. A typical massage of this sort might look likes

. OPEN 15

Such a command always starts with symbol # and ands with either a linefeed or a rubout, depending on whether the user is satisfied with the command or wishes to abort it. The only exception to this rule is the specific command

.

which inserts an # in the data atreem to the Host. Commends may occur anywhere, and need not start on a new line. Upper and lower case may be freely intermixed in the commend.

Between the # and the linefeed there will typically be one or more words to identify the command; perhaps followed by a single perameter. The TIP is not very sophisticated; and thinks the only important thing about a word is its first letter. This permits the user to abbreviate a bit; the more usual rendering of the first example might be:

#0 15

Once the user has started typing the parameter of a gommand the old value of the command will have been destroyed, and dannot be recovered by eborting the command.

Almost without exception the effect of a TIP command is to set a parameter or made for the terminal. Even apparently direct commands like

P OPEN 15

(which initiates an elaborate exchange of messages resulting in a connection to the remote Host system) actually set a mode flag to request the appropriate action when the TIP is free to undertake it. To understand the TIP behavior is really to understand the complete set of parameters and the commands to change them. Normally, any parameter can be changed at any time by the user at his terminal. Exceptions occur when the user tries to change connection parameters on an open connection. An #OPEN 13 executed while talking to Host 15 would generate the error message "Can*t" (the connection to Host 15 must be closed before a connection can be opened to Host 13).

Commands often consist of several command words; for example,

. DEVICE CODE ASCII

Such commands may be abbreviated; for example

P D C A

The spaces are required; # DCA is not a legal command, Upper and lower case letters may be freely intermixed.

An unusual variation in command format is to place a number between the # and the first word of the command. In this quae, the command is not meant for the terminal typing but for the terminal attached to the port of that number on the same TIP as the user.

In the normal course of things, a user will go through five more or less distinct staces in typing into the net. First, he will be concerned with hardware-power, disling in, atc. Then he will establish a dislogue with the TIP to get a comfortable set of parameters for this usage. Next, he will instruct the TIP to open a connection to a remote Host; and finally, he will mostly ignore the TIP as he talks to the remote Host. The following sections will describe these stages in more detail.

A. Hardware Stage

The hardware stage is out of the scope of this section, except for the final stage of connection. When the TIP detects a terminal on one of its previously idle lines it normally goes

into a "hunt" mode. In this mode it expects the very first character it sees to describe the terminal, according to the following scheme

> ASCII Terminals at 110, 150, or 300 baud type E. (Note that this must be upper case.)

2741 Correspondence Terminals type], 4, o, or 1 depending on the element used with the terminal -- see Table 4-A,

274: PTTC Terminals type:

- 6 for model* 938, 939, 961, 962, or 997
- o for model 942 or 943
- w for model 947 or 948
- 1 for model 963, 996, or 998

ASCII Terminals transmitting at 110 but receiving at 1209 baud type D. (Again, upper case)

The TIP will deduce terminal rate, character size, and code conversion based on the character typed. When the TIP makes its dealsion it types out TIP's name in the terminal's own language followed by the version number of the TIP software system and the octal post number. Then it is ready to go,

Establishing Parameters 8.

In stage two, the user is denomined with initializing parameters. The naive user should skip stage two and accept the TIP's default parameters until an obvious problem arises. The following questions are answered in stage two:

- is When shall the TIP send off messages to the remote Host?# Here there are several options. (The TIP is initialized to send on every characters which is simple but inefficient.)
 - TRANSMIT NOW
 - . TRANSMIT ON MESSAGE-END
 - # TRANSMIT ON LINEFEED
 - . TRANSMIT EVERY N

 with reset the TIP to its initial state, transmitting every character. If the parameter to TRANSMIT EVERY is a large number (e.g., 250) the TIP will seve up as many characters as it can before sending a message, but does not offer any quarantee that the total number specified can be buffered.

2. Who shall egho, and when?# Eghoing is a complex problem, without any nest solution. We have chosen to give the user the means to tell us how he wants it done, since it is hard to guess correctly in advance. Basically, eghoing gan occur at the terminal hardware, in the TIP, or in the remote Host. The corresponding TIP commands are:

ECHO HALFDUPLEX (Echo at terminal)

ECHO ALL (Eaho at TIP)

ECHO NONE (Echo at remote Host)

ECHO REMOTE (Send TELNET "remote echo" and

perform internal ##E#N)

ECHO LOCAL (Send TELNET "local acho" and

perform internal ##E#A)

In the ECHO NONE mode, although shapedters for the remote Host are not echoed, the TIP will echo commands.

Network protocol specifies that echoing shall start out

in the *#ECHO ALL or *#ECHO HALF modes. The TIP will try to guess from the terminal type which of the two is appropriate. The goal of edhoing strategy is to avoid the unreadable alternatives of the blank page and the doubling of every character. The naive user is advised to accept the TIP's default parameters until trouble of this sort arises. A 2741 is incapable of changing echo mode; it is always echo halfduplex.

To allow more complex echoing conventions, the TELNET protocol provides a mechanism whereby the remote terminal user may instruct the serving Host whether or not to echo characters. The ECHO REMOTE and ECHO LOCAL commands at the TIP allow TIP users to use this mechanism after the connection is made.

Finelly, many Hosts which provide service request the TIP to allow them to do the etholog. The TIP elways grants this request (even for 2741 terminals). The user, if he does not desire this mode, must cancel it AFTER the connection to the Host is established.

C. Connection to Remote Sites

In stage three, the user is concerned with establishing a connection to a remote site.

P OPEN 15

This amounts to "met the Host number parameter" and "edd the user to the queue of users welting for the TIP's connection mechanism", Appendix A lists the Host numbers of all the sites querently in the network,

Connecting to a Host requires establishing a bimdirectional link, so that the terminal can send cherecters to the Host and vice versa. The request to connect to a Host is thus really a request to establish both transmit and regards sections.

When the user reaches the head of the queue weiting for the TIP's "connection" mechanism, the TIP will type "Trying...".

Following the message "Trying", the user will receive some of the following messages:

Open success

Not Trouble remote site cannot be reached

Refused remote site up but refusing

Host Scheduled Down Until Set, at 1850 GMT

Host will be back up at time and date

indicated

Host Unscheduled Down Until Sat, at 1858 GMT

Host will back up at time and date

indicated

Host not responding

Remote site not up, unknown when up service will resume

ICP Interfered With

The Host has not performed the ICP correctly and the TIP has refused to open a connection.

The connection mechanism will run continuously until a state described above occurs. This can be annoying when the remote site is obviously not going to respond. The command

● CLOSE

will abort the gurrent connection attempt. The user is then free to reattempt to open the connection or to attempt to open a connection to some other Host.

The TIP's connection mechanism has caused users some problems. Perhaps a discussion of what the connection mechanism is doing and how it works will alleviate some of the grief.

First of all, users attempting to connect to different Hosts will never interfere with each other, although users simultaneously attempting to connect to the same Host will be serviced serially.

For the user, opening proceeds in three phases. In the first, the user is queued up waiting to "get" the TIP's connection mechanism. In the second, the user has gotten the TIP's connection mechanism and is beginning the sonnection sequence. In the third, the user has completed the connection sequence and is waiting for the Host to open up the actual data connections. Henry of the problems stem from the fact that only one user may be proceeding through phase 2 at a given time to a given Host, Hense the the TIP types out "Trying" when you get off the queue and the sonnection mechanism begins trying to open your connections. Thus the "Trying" message signifies the trensition from phase 1 to phase 2.

D. Use of Remote Sites

In stage four, the user is normally talking to the Host without concern for the TIP, All the TIP gommands are still evaluable.

One command that will eventually be of interest here is

· CLOSE

This command starts the shut-down procedures. The TIP will echo "Closed" when the process is finished. The TIP does not know how to log you out of the remote Hosts. You must do this yourself before closing the connection.

The virtual terminal has a key labeled "BREAK", Some real terminals have a break key, and some Host systems expect to see breaks. Those terminals with a break key (but not the 2741 ATTN key) may simply use it. Others must type the command

. SEND BREAK

The interpretation of the break is entirely up to the reseiving Host == meny Hosts (gnore it.

The virtual terminal also has a key labeled "SYNC". No real terminals have such a key, and the function is unique to network use. The "SYNC" key is a clue to the remote Host that there is an important message which seems to be buffered in an "inscressible" place. The TIP and the Host go to some trouble to get the SYNC indication over a different channel which bypasses the normal buffering sonventions. The command to send a SYNC is

. SEND SYNC

Typical usage of these commands might be # \$ 8 followed by # \$ 8.

As stated earlier, the TIP nominally treats a carriagereturn typed by a user as a derriage-return/linefeed. The user
may gause the TIP to treat derriage-return as only
carriage-return by executing the command

OCLEAR INSERT LINEFEED

To return to the nominal mode of cerriage=return/linefeed, the command

PINSERT LINEFEED

should be executed.

If at any given time the user types characters faster than a server Host will take them from the TIP, the TIP discerds characters it can not buffer and echos them with an ASCII BEL (on a 2741, the type element is wiggled).

The user may sometimes use a server Host with which it is desirable not to have # be a TIP reserved character. The user can change the character which introduces TIP commands using the command

PINTERCEPT N

By typing #INTERCEPT followed by a decimel number representing an ABCII character, the user changes the TIP command character for his device to the ASCII character represented by the number. The INTERCEPT ESC command resets the TIP command character to at-sign (#). Thus,

SINTERCEPT 42

*INTERCEPT ESC

changes the TIP command character to asterisk (*) and beck to atwaign (*) assuming the device was in the nominal mode (*) before the first command was executed:

If the user attempts to change the intercept character but feils to type a valid decimal number (or a character string beginning with E) the TIP will type the diagnostic "Num" and will set the intercept character to at=sign.

E. Connection Loss and Restoration

Starting with Hosts running TENEX Version 1,32, if TENEX helts, the TIP will notify users connected to it of this fact by typing "Connection suspended". At this point the users are free to do one of two things. First, they can welt till TENEX restores service, in which dose the TIP will type out "Connection Restored" for if after the the service interruption the connection could not be restored, the TIP will type out "Host has reset connection"). Alternatively, the user is free to open a ponnection to any other Host, in which case the TIP will invisibly glose the TENEX connection. It is also important to point out that if a user lust leaves his terminal unattended ecross a TENEX service outsign without releasing the connection (any network related command such

as PH, PO, PN, PC will do the Job) his Job, directory, etc., are left at the mercy of envone who acquires that terminal.

F. TIP News and User Feedback

There is frequently information which the group developing and debugging the TIP system wishes to donvey to TIP users. instance, when a bug is detected, we may wish to warn users not to use a sertain feature until the bug is fixed. When a miner improvement is made, we may wish to notify users. Further, there is frequently news about the state of the network or the state of a particular Host which should be conveyed to TIP users. Finally, TIP users may wish to communicate with the TIP development staff or the Network Control Center staff about problems or suggested improvements for the TIP or the network. Consequently, we have constructed a mechanism which we hope will provide for communication in all the above directions. This mechanism is the Network Virtual TIP Executive a To activate this mechanism, the TIP user may give the TIP command #N. (The same one which handles the user's login to the TIP:) This commend causes the TIP to perform the negessary protogo; to make a gonnection to the Network Virtual TIP Executive which posides on several of the network TENEX systems. Once the Network Virtual TIP Executive activated, we think its operation is self-explanatory, Presently available features within the Network Virtual TIP Executive are a Network News feature, a Host Status feature, and a "Gripe" feature. The latter provides users with a mechanism for sending messages to the TIP development or NCC staffs. We recommend that TIP users get the network news at the beginning of every TIP session (during the login prodedure).

When a user issues an eN command, the TIP requests support from ALL cooperating servers. Thus, the user should be able to reach a news facility, somewhere, elmost all of the time. However, in the event that no cooperating server is available the TIP will time out the eN command in about thirty seconds (and respond with "Timeout OK to proceed"). An eC command will about an eN immediately.

Of course, TIP users with an immediate need for communication with the NCC or TIP development staffs should telephone (collect) the Network Control Center (617-661-8188). Users with general questions about network usage (How do I find out if Host X is ever going to be up again? What's happening with a Host/Host protocol for graphics?) may also call the NCC.

ARPANET Completion Report Chapter III

DRAFT Section 2 Design and Implementation

Table 2-4. TIP MESSAGES TO THE TERMINAL USER

BAD The TIP doesn't recognize the command

Closed Connection closed, usually by server Host

Connection Restored

Destination Host has restored the connection as it was before the Host helted.

Connection Suspended

Destination Host has haited operation.

Host broke the connection

The Destination Host's service is restored but all network connection tables have been reset.

Host not responding

Destination Host not up from the network's point of View. It is not known when service will resume.

Host Scheduled Down Until ...

Destination Host is acheduled down until the date and time indigated.

Host Unscheduled Down Until ...

Destination Host is unscheduled down until the date and time indicated.

ICP Interferred With

The Host has not performed the ICP correctly and the TIP has refused to open a connection.

Login first

The year must login to the TIP before using this commend.

Net Trouble Destination IMP gannot be reaghed due to some kind of trouble in the network.

NO Parameters may not be set for specified terminel.

Num The TIP expected a number am gommand terminated,

Open Connection opened by server Host.

Section 2 Design and Implementation

Q Refers to the Reseive side of a connection.

Rafused The remote Host rejected the attempt to establish

connections.

Ť Refers to the Transmit side of a connection.

Timeout OK to proseed

No RSEXEC is available to escept a login = the user mey continue this TIP session without

accounting or access control.

TIP GOING DOWN The TIP is going dawn in the number of minutes

indicated -- Quickly stop what you are daing

and stop using the TIP.

TIP NAME

The TIP heard the user dial in and establish rate. The number following NAME is the TIP software system version number. It is followed by the

octel part number.

The TIP is now servicing the user's OPEN request. Trying

Wait The TIP is ettempting to contact an RSEXEC for

user todin.

DRAFT dection 2 Design and Implementation

Table 2=5. TIP COMMAND SUMMARY

BINARY INPUT END
Leave 8-bit binary input mode

BINARY INPUT START
Enter Sebit binary input mode

BINARY OUTPUT END
Leave Smbit binary output mode

BINARY DUTPUT START
Enter 8-bit binary output mode

CLEAR DEVICE WILD

Set device to be unwild

CLEAR INSERT LINEFEED
Stop inserting iinefeed efter cerriagemeaturn

CLOSE
Close all outstanding connections, or abort current Host tegin

DEVICE CODE 37
Establish parity semputation for Model 37 Teletype

DEVICE CODE ASCII

Establish code conversion for an ASCII terminal

DEVICE CODE EXTRASPADDING Establish code conversion for a terminal with slow CR

DEVICE CODE OTHER-PADDING

Establish code conversion for a line printer

DEVICE RATE _#
_# is a 10-bit code specifying hardware rate and
character size settings

_# DIVERT OUTPUT

Capture device _# and divert this terminal*s output

to it, _# is an octal number.

ECHO ALL
Local TIP=generated echo == TIP eghoes everything

ARPANET Completion Report Chapter III

DRAFT Section 2 Design and Implementation

ECHO HALFDUPLEX
Terminal=generated eaho == TIP echoes nothing

ECHO LOCAL

Send the Telnet "ECHO LOCAL" pherester and perform internal E A

ECHO NONE

Remote Host-generated echo for data ==
TIP echoes commands

ECHO REMOTE Send the Teinet "ECHO REMOTE" character and perform internal E N

FLUSH
Delete all characters in input buffer

_# GIVE BACK
Release control of deptured device _#;
_# is an octal number.

HOST _# Simultaneous FPS T HF and FPR F HF

INITIAL CONNECTION PROTOCOL

Start the initial connection protocol

INSERT LINEFEED
Insert linefeed after dereiage=returns

INTERCEPT _# Use _# as TIP command character

INTERCEPT ESC Leave 7-bit binary mode

INTERCEPT NONE Enter 7-bit binery mode

LOGINALA An obsolete form of OPEN

M _# _# Mag tape dommand _# with argument _#

ORAFT Section 2 Design and Implementation

NETWORK=VIRTUAL=TIP=EXECUTIVE

Connects the user to the NetworksVirtuals
TIPSExecutive, also used to initiate login for some terminals.

OPEN _M

Open a biedirectional connection to the Host decimal address is specified

PROTOCOL BOTH

Simultaneous #0P T TH and #0P T R#

PROTOCOL TO RECEIVE

Manually initiate connection protogol

PROTOCOL TO TRANSMIT

Manually initiate connection protocol

RECEIVE FROM HOST _#

Establish Host _# parameter for manual initialization

RECEIVE FROM SOCKET _#

Zetablish masket._# parameter for manual
initialization of connection == socket _# is
often in octai

RECEIVE FROM WILD

Equivalent to ##R P S dany>#

RESET

Logout surrent fir connection and hang-up data set

Reset NCP

Resets NCP

SEND BREAK

Send the Telnet "BREAK" charmater

SEND COMMAND

Send the sommend escape thereater

SEND SYNC

Send the Telnet "SYNC" sharester and an "INTERRUPT SENDER" message

SEND TO HOST _#

Establish Host _N parameter for manual initialization of connection

DRAFT

Section 2 Design and Implementation

SEND TO SOCKET _#

Establish mocket _# paremeter for manual initialization of connection == mocket _# (a given in octal)

SEND TO WILD
Equivalent to "#5 T 8 <anv>#

SET DEVICE WILD

Equilibrium to the commands "FR F H damys",

#85 T H damys", "85 T S damys", and "88 F S damys".

TRANSMIT EVERY __#
Send off input buffer at least every _#th
character where E4_#4256

TRANSMIT NOW Send off input buffer now

TRANSMIT ON LINEFEED

Send input buffer every time a linefeed is endountered

TRANSMIT ON MESSAGE-END

Send input builer every time an end-of-message
is encountered

2.3.4 The NCC Hardware and Software

2,3,4,1 Introduction

The NCC Host has two general jobs; it is a tool for diagnosing network problems and it is a network data tollector.

As a tool for diagnosing network problems, the primary NCC job is to use the periodic IMP messages (plus the internal clock timeouts) to discover changes. If the changes are major, the NCC Host will sound an elerm end flash a light, to incure that the change is noticed quickly. All ghanges, whather major or minor, are presented as a human-engineered printout, called the LOG. The LOG reports not only complete fellures of a line or an IMP, but also information about unusual events which often precede complete fellure (e.g., a noisy line may report a low error rate as some hardware component weers out). The IMP programs also detect many conditions which are not normal, but recoverable, and sends special "trap" messages that appear on the LOG.

The NCC Host is evere of the network topology. Unlike each IMP, which knows only how many "hope" are necessary to reach enother IMP, the NCC Host reports about "lines" to the network support staff. Each IMP sends the NCC a snapshot of its environment. The NCC integrates the anapshots into a single picture of the network as a whole, and presents the larger view

with lights, elerms, and a LOG printout, This knowledge of topology is occasionally of great importance; for instance, if the network ever partitions, isolating several sites, the NCC gen compute the "frentier" and flash those lights while turning off the lights for the "invisible" sites and lines,

As a network date dollector, the NCC Host saves some of the raw IMP data in some memory, Host and line throughput data is collected from each IMP and summarized on a terminal hourly. The status of each IMP and line is saved every 15 minutes; a status summary printout is produced every 8 hours. The compected encoding of the LOG printout builds up in memory. This data may be shipped to other Hosts upon request.

The data collection abilities work in gonjunction with processes running on Tenex (PDP=10) somputers at other network Host sites. The full data collection is a distributed computation. Other network Hosts have reliable bulk storage and general-purpose languages — the NGC Host is a mini-computer without much storage or power. The two kinds of machines work together to get interesting data to the proper place for analysis. Typically, the NGC data is used for "batch" rather than "real-time" computation on other Hosts.

2.3.4.2 The NGC Hardware

The NCC Host looks like a Terminal IMP (TIP). It is a Honeywell 316 with a Multi-Line Controller (MLC) added. It has 28K words of gore memory, an IMP interfece, a light-display/alarm, and several terminals. It has a paper tape reader and a modem interface for program loading. It does not need apecial instructions (for instance, multiply and divide), disks, or other desirable options; the hardware is almost identical to the BBN=TESTIP. In case of serious hardware failure, the BBN=TESTIP may be substituted for the NCC Host with a minimum of effort. The

- 1) IMP Interface The NCC's IMP (nterface is identical to the IMP's Host interface; the interfaces are connected by a cable which juggles the IMP and Host definitions appropriately. All network data, whether IMP or Host traific, uses this path.
- 2) Primary terminal a This Model 35 Taletype has two major functions. It prints periodic reports of Host and line throughput and status, Segondly, it is used to change NCC parameters, and as a debugging sid, much as the IMP Teletype and DDT gods are used for IMP control and debugging.

- 3) Lightbox & Alarm = A SSNedesigned I/O device with 64 LED lights, i6 gwitches, and an audible buzzer; this is the primery werning of important network changes. Since there is only one such device, the progrem can run without it, although the support staff would have to keep a careful eye on the LOG printout if the lightbox/elarm were broken.
- 4) LOG terminals The LOG produges a hard dopy record of network events. Normally two small lineprinters are used:

 a Centronics 306 is next to the NCC Host, and a quiet Inktronics is in the manned area of the Network Control Center; both run at 1200 baud.
- 5) MSG=DUMP terminal a One terminal is dedicated to printing an outal dump of any network message which domes from an IMP (as apposed to a Host) but is not a reaugnized status or throughput message. This facility is a flexible IMP debugging aid, particularly for viewing "bad" peckets.
- 6) TTHEXER terminals = Other MLC ports may be used to exercise terminals by continuously printing a simple pattern. This is an aid to hardware maintenance of the LOG and MSG=DUMP devices.

2.3.4.3 NCC Outputs

The following few pages contain sample printouts. The first page shows the LOG printout similar to the lineprinter output in the Network Control Center. The remaining pages of printout dome from the primary TTY, often referred to as the "summery TTY". The IMP and line status printouts are compressed, with the following 9 single characters used to indicate different status:

- I Unknown
- Down in the plus direction (lines only)
- Down in the minus direction (lines only)
- * Down
- ع نا 🐞
- Plue and looped (lines only)
- Minus and looped (lines only)
- Z Both ends looped (lines only)
- X Status changed one or more times during the interval

The Host throughput data is divided into 5 gategories as follows:

Mean Messages to Hosts at other nodes,
Men Messages from Hosts at other nodes,
Pean Packets to Hosts at other nodes,
Pean Packets from Hosts at other nodes,
Meas Messages to Hosts at this node,
Peas Packets to Hosts at this node,
Peas Packets to Hosts at this node,
Peas Packets to Hosts at this node,

Every 8 hours, the IMP status, line status, Host throughput, and line throughput reports are printed, summerizing the period state midnight. Every hour, the Host and line throughput reports summerize the previous hour. The "quisk summery" and "host summery" reports print only at the operator's command.

```
NOVEMBER 13 1976 ARPA NETWORK LOG PAGE 99
1100
1100 IMP 601 *(102.,4,0)
     LINE 411 ERRORS PLUS 2/16
      IMP 58: HOST : DOWN
1101 LINE 411 ERRORS MINUS 1/17
      IMP 611 +(102, 0,0)
      IMP 581 HOST 1 UP, HOST 1 DOWN
      IMP 608 = (102.,1,177771)
      IMP 581 HOST 1 UP
1102
     LINE 412 ERRORS MINUS 1/17
      IMP 618 *(102..0.0)
      IMP 561 HOST 1 DOWN
1103
      IMP 60; +(102,,177775,0)
      IMP 58: HOST 1 UP
      IMP 581 HOST 1 DOWN
1104
      IMP 451 HOST 1 DOWN
      IMP SEE HOST 1 UP
      IMP 453 HOST 1 UP
      IMP 611 *(102..0.0)
1105 LINE 411 DOWN ***
      IMP 58; HOST 1 DOWN, HOST 1 UP, HOST 1 DOWN
      IMP 581 HOST 1 UP
1186
      IMP 611 *(182..0,0)
      IMP 581 HOST 1 DOWN
      IMP 60: +(102,,177777,177776)
1107
      IMP 38: HOST 1 UP
      IMP 58; HOST 1 DOWN
      IMP 39: TRAP: (7786,2,1)
1108
          61 HOST 0 UP
      IMP
      IMP 111 + (32, 104647, 0)
      IMP 111 HOST & DOWN
1189
1115 LINE 41; ERRORS PLUS 1/16
      IMP 411 + (63.,400,47)
1116 LINE 411 UP ***
      IMP 30: POWER FAILURE, RESTARTED
1117
     IMP 30: OVERRIDE ON, TRAPE (23456,0,36)
1118
     LINE 881 LODPED PLUS ***
```

NOVEMBER 4 1976 ARPA NETWORK SUMMARY 0000+0800 PAGE 148

TIME		-	1234567890		•	F - ',
8688		********	*********	*********		*********
8838	. 7 .	********			X * * * * * * * * *	
0045	11111111			*********	********	********
0100		********	******	********	X	*******
0115	********	**********	*******		*****	*****
0130			*******			
0145	**; **; **;			******	X.,,,,,,,	
8288	77.		*******	111111111	*******	,,,,,,,,,,
0300		********			X * * * * * * * *	********
0315			*****			
0330						
8345		•••••		*******	X	*****
9499						
0500	111111111			********	X	
0515	11011111	*********	*******		********	********
0530	*********	********	*********		X	*******
9545				********	*******	*******

NOVEMBER 4 1976 ARPA NETWORK SUMMARY 8088800 PAGE 150

LINE STATUS

TIME	1234567890	1234567890	1234567890	1234567890	1234567890	1234567890
9999	******		******	.,,,,,,,,	*******	,7,,,*,,,
8938	1 ⁷ , 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,	**************************************			X	• • • • • • • • • • • • • • • • • • • •
9945	17,0,,,,,,		*****	********	*X	,?,,,*,,,
0100	7,0,0,0	**************************************			XX******	• • • • • • • • • • • • • • • • • • • •
0115	7,,,,,,,,	?			******	
0130	11,1,4,1,1	111111111			xx******	
9568	17,1,1,1,1	111111111			********	
0300	17.11.11.1	111111111			XX	
0315	********	********			xx******	
9330		********			********	
8345	.7,,	****			XX*****X	
8400	. 7	********			•••••	
8500	********	111111111			X	
0515	17, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,	7,7,,44,			X	
0530	*********	***********			X	
0545	17,10,000				X * * * * * * * * *	
0600			******	*********	x	171111****

NOVEMBER 4 1976	ARPA NETWORK SUMMARY	0000-0800 PAGE 152
H087 0	H087 1 H087 2	HOST 3 THROUGHPUT
SITE 1 UCLA MEDN 23778 MEDN 26460 PEDN 23778 PEDN 27937 MEDS 4667 MEDS 1677 PEDS 4874 PESS 2862	7383 2489 14963 3284 12117 2489 10976 3284 1152 1757 1137 4562 1152 2742 1344 4562	
SITE 2 \$RI=5 MmbN 61632 M<=N 67072 P=bN 61632 P P P M=b8 2676 M<=3 4311 P=b8 2677 	125480 124160 125696 132416 5260 3763 6810 3814	
SITE 3 NUC MODN MODN MODN PODN PODN	539 721 539 956	
SITE 4 UTAH Medan Medan Pedan Pedan Pedan	52928 36408 52928 38336	
SITE 5 88N M=>N 134400 M<=N 172160 P=>N 148672 P<=N 177856 M=>8 129600 M<=\$ 177344 P=>\$ 143168 P<=\$ 177792	365836 18 335486 385636 16 350016 513 217152 481 147904 513 217472 2036 162816	13979 13911 21815 14062 31506 53056 34624 53120

NOVEMBER 4 1976 ARPA NETWORK SUMMARY 1680-1709 PAGE 177

LINE	TH	ROUGHPI	JŤ.						
	12345678; IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII	PHMUUUUUSBSASHMULCCH8MCNAIGOM	10000000000000000000000000000000000000	N T X SRTP N I AB JOR ENTR NOV SAE STOCCT JAPCOCKREJIIIGT LMGJUERMANKGAAC GCROSUCKBIRSCHLOREYMROSMIRODC EOWDUNSXALSINMIRABNAARNAKASS	24430176895376427266270520324 6124377689537683854057689537683168386768316838676837683868386838687683 2436768953168388568768384798888888888888888888888888888888888	N T X JATP N I AB JOR = N=1 EOV SAB STOCK GCROSUCRBLIFICT LMGJUSPMAERGAAC ECWOUNSXALSINMIRABNAARNAKASS		A AAAB BT CHANNAN B 2 CR MUUUUUSBBASHMUULGGHSMCNAIGDM	846291948948722172862488 7116291948988483797172862488 71198918988483797172862489 711036736893797172862488 711114112848722172862488 711114112848722172862488 71114112848722172862488 71114112848722172862488
LINE LINE LINE LINE LINE LINE	_	MITRE BBN T BBN GWC PENT NBS XERQX	10000000000000000000000000000000000000				TO TO TO TO TO TO	MITRET BBNGT DD PBRNDC PBRNDC RANNC RANNC RAMES	

1502 NOVEMBER 10 1976 GUICK SUMMARY

IMP

17.
1234567890 12

IMP 1 MEM PROTECT OFF

IMP 2 MEM PROTECT OFF

IMP 3 MEM PROTECT OFF

IMP 4 TIP UP

IMP 5 MEM PROTECT OFF OVERRIDE ENABLED

IMP & MEM PROTECT OFF

IMP 7 TIP UP

IMP 10 MEM PROTECT OFF

IMP 11 MEM PROTECT OFF

IMP 12 MEM PROTECT OFF

IMP 13 TIP UP

IMP 14 MEM PROTECT OFF

IMP 15 MEM PROTECT OFF

IMP 16 TIP UP

IMP 17 TIP UP

IMP 18 TIP UP

Section 2 Design and Implementation

1507 NOVEMBER 10 1976 HOST SUMMARY

IMP

1234567890 1234567890 1234567890 1234567890 1234567890 1234567890 123

IMP 1 HOST 0 UP HOST 1 UP HOST 2 DOWN

IMP 2 HOST 0 UP HOST 1 UP HOST 2 DOWN

IMP 3 HOST & DOWN HOST 1 DOWN HOST 3 UP

IMP 4 HOST & DOWN

IMP 5 HOST Ø UP
HOST 1 TARDY
HOST 2 UP
HOST 3 UP

THP 6 HOST 3 UP HOST 1 UP HOST 3 UP

IMP 7 HOST Ø UP TIP UP HOST 3 UP

IMP 9 HOST 0 UP HOST 1 TARDY HOST 3 DOWN

IMP 10 HOST 0 UP
HOST 1 DOWN
HOST 2 DOWN
HOST 3 DOWN

2,3,4,4 The NCC Software

This section describes the NCC progrem organization; when reasonable, reference to specific subroutine names in the code of the NCC system has been made.

The progrem is structured in three logical levels a background, foreground, and I/O. Background, the idle machine state, is used as a machanism for starting up some I/O events. Foreground processing, the bulk of the program, includes handling IMP messages. I/O routines are short, high-priority, unwinterruptable strips of the program.

The Honeywell 316 has an interrupt atructure which permits multiple interrupt levels. The hardware instruction set allows both plobel and selective inhibit/enable control. The background attate means running with most or ell interrupts anabled, but not eurrently on any interrupt levels. The foreground state means operating on either of the two lowest priority interrupt levels, with these two levels disabled, but with other interrupts allowed, I/O interrupts num to completion with the global inhibit in effect. Different hardware levels must work together for many of the NCC functions. Some of the coftware interfaces between hardware levels are very simple, yet some are complex. Most I/O interrupts signed the main IASK processing in some way. The interrupts signed the main IASK processing in some way. The

IMP interface input and TASK, IMP interface output, Background, and TASK, MLC outin (LOGGER), MLC output, Background, TASK, and CLOCK, Primery TTY, Background, and CLOCK, TASK and CLOCK.

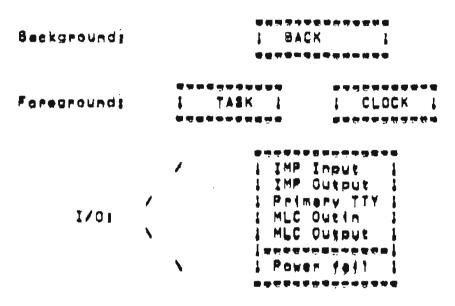


Figure 2-18: The Three Logical Levels

Foreground processing deause on two lowerfority interrupt levels delied TASK and CLOCK, GLOCK, delied by a 25.6 millisegond hardware interrupt, keeps track of the time of day, notices if any IMP or line has not been heard from for several minutes, and either performs or requests background to perform time-dependent functions. GLOCK is covered in more detail later, GLOCK and TASK do not interrupt each other, and can share subroutines.

The TASK interpunt is a "software" interpupt coused by execution of an I/O instruction. The IMP interface input causes the TASK interrupt whenever a full message arrives from the Figure 2-11 is the TASK overview, showing the four major sections of code within the wavywedged boxes. Newly arrived network messages are routed to one of four places a IMP status handling, IMP throughput message handling, IMP octal printout (M8G=DUMP) handling, and Most core dump request handling, The status message handling was govered in the "IMP MESSAGE PROCESSING" chapter. The other three message handling routines are short and straightforward; the throughput values are added into the proper Host and line tables; the MSG=DUMP message is copied to another area for later printout; and the core dump messege shocked for legitimacy, then the parameters are saved in the NONPWC/NONPLO variables.

Figure 2-11: Major Portions of TASK Level

CLOCK is the second interrupt which is used in the foreground level of the program. It keeps track of time, notices IMPs and lines which have not reported to the NCC recently, sets flags for background level to initiate periodic printeuts, and initiates periodic disk file updating. If the optional alarm and light display hardware is installed, it is updated frequently,

Since the clock can be locked out far in excess of its 25,6 milliaecond period, the main entry, TOR, computes the number of intervals which have occurred since the last clock interrupt and calls TOWORK once for each elapsed interval. As explained in the section about program structure, I/O interrupts are enabled but TASK and extra CLOCK interrupts are inhibited.

Every Second

If the primary terminel timer is running (and inhibiting the output-only function), it is counted down. If sero is reached, the primary terminal state is set to "idle" (TT88Y=0),

If senseswitch 3 has been turned on, GLOCK checks for an inective machine. The IMP+Host ready line is dropped. The LOG terminal must be idle and the primary terminal must be idle. If everything is inective, the program halts.

Every Hour

The primary terminal hourly throughput printout flag is set for action by background level. If the hour is a multiple of eight, the total throughput and total status printout flags are set. If the hour is 24, the LOG terminal date and page number are est up for a new day.

If the "symmery" time has reached the hour 1 or 25, the tables OHTP and OLTP are sleared.

The tables HTP1 and LTP1 are copied into HTP2 and LTP2. The HTP1 and LTP1 values are added into DHTP and DLTP4. The HTP1 and LTP1 tables are seroed.

Every 15 Minutes

The IMP and line status values are pushed on a stack that sontains 96 such entries = one for each 15=minute segment in a day. In this way, status attes are quantized to 15=minute periods, and any information on atatus since the previous midnight can be recovered.

Every Minute

The program checks IMP and line Status each minute, If an IMP has not reported resently, the following actions are takens

- 1) The IMP is deglared down in the IMP status table
- 2) A logger message is printed to indicate the IMP is down
- 3) The eterm display light for the IMP is flashed
- 4) All previous status indicators for that IMP are pleared to nominal values

Na In∮o Up Down, so Down with Looped | Unknown | Up | Unknown | Unknown | Looped | Na Infa 7 1 , 1 7 1 M Up Up : Up i 1 (mbo i 1 fmbo i 1 fmbo 1 N | Unknown | 14mbs | Down | Down | Looped | erfor# 9100 Down with lunknown | 1(mbo | Down | Down | Looped | minus 1 N l plus D Looped i Looped i limbo i Looped i Looped i Looped i minua | l minus i minus i both

PLUS END

Figure 2-12 Line Status Determination

(Plus and # highenumbered IMP)

Line status is determined every minute by comparing the latest report from each of two IMPs on a line. The transition table shown in Figure 2-12 is used to choose the current state.

a change in status is detected, the following actions are takeni

- 1) A logger message is printed
- 2) The new status is recorded in the line status table
- The alarm display lights are set 3)
 - a) UP turn on both plus and minus lights, turn off flashing
 - DOWN MINUS = turn on plus, turn off minus lights, turn on flashing
 - DOWN PLUS turn off plus, turn on minus lights, e) turn on flashing
 - ALL OTHER STATES # turn off lights, turn on flashing d)

Each minute, IMP and line statuess are updated, IMP visibility is a central concept to this procedure. There are two possible reasons why an IMP might not be reporting int because the IMP is actually malfunctioning, and second, because it has become implated away from the NCC (i.e., finvisible), The NCC distinguishes these games so that in the event of a failure which partitions the network, the NCC can determine the IMP involved, thus minimising the number of irrelevant Logger messages and making such occurrences less of a deduction problem for the operator.

2.3.5 The Pluribus IMP Hardware and Software

In (te original form, the IMP was based as much as possible on the most suitable offethershelf hardware then available; apecial hardware design was kept to a minimum and consisted of interfaces and a few special features which were added to a standard machine. During the first few years of the networkfs existence, new and more flexible computer structures began to appear on the market, and the special requirements of packet=switching began to be better understood. The Pluribus architecture, developed apacifically to suit the needs of a packet-switch, is the outgrowth of these changes. There are two primary goals for the maghine, representing areas in which the earlier IMPs were felt to be lacking. These were flexibility ().e., the ability to expand or contract amouthly over wide ranges) and reliability. Originally there was a more primitive goal of higher throughput, but this was seen seen to be balanced by the need for a cheaper, smaller maghine with less throughput, Bendwidth is thus seen as one domain in which greater flexibility is desired.

Let us consider the (saue of flexibility in a little more depth. In most mechines dertain handware futilities, are shared smong the various logical units. These include reak apage, power, cooling, etc. Generally these utilities come in fairly large

chunks, with correspondingly large steps in cost. Thus, one can typically add, say, interfaces up to a certain point; at which time a new reck, power supply, etc., must be added to permit further expansion. Even then, one may run out of logical channels or come up against other hard boundaries. The 516/316, for example, has a fixed memory channel errangement which limits connection to a total of at most seven high-speed circuits and/or The specific component which is hest computers. inflexible in most systems is the processor; that is, there typically no precessor modularity or possible varietion of processor capacity. The flex(bi)(ty coel of the Pluribue smooth large step functions in east by utilizing a highly moduler design and to bush really hard boundaries (sugh as absolute limits on memory addressing especifities or processing especity) well beyond requirements anticipated at least for the next few years,

The machine design was thus to allow for large numbers of I/O units, for wide ranges in processor power (up to at least an order of magnitude in traffic bandwidth handling capability improvement over the 316), and for larger possible memory (to permit longer and/or faster links, say, via setallite).

Now consider reliability, If a single IMP fells on an average of once a month (quite a high MTBF for conventions) computers), then in a network of thirty such IMPs, one will fell

on an average of once a dey. In a network consisting of large numbers of computers (e.g., hundreds), issues of reliability take on paramount importance. In a multiply connected network the failure of a single IMP should be felt only locally. Equnomics, however, limits the number of paths between any given pair of IMPs. In the ARPANET meny IMP pairs are connected by only two disjoint paths, some by only one, and this tends to increase reliance on individual IMPs and to make their reliability a key issue.

If a single computer fails on an average of ten times a year, then a collection of ten computers, treated as a unit, will fail on an average of 100 times a year. However, suppose that rather than viewing the ten computers as a unit which is down if any one of its constituent computers is down; we view the ten computers as a unit which is up as long as any of its constituent computers. Is up. Further, suppose the mean time to repair a failed computer (a small compared to the time between failures. In that case the probability that the constituent computers will all be down is very small, so it is unlikely that the unit as a whole will be down. The reliability of the Pluribus IMP takes advantage of such probabilities. Note that a key assumption is that individual are independent of one enother. Although it is impossible to guarantee this independence (flood; total power failure, mabotage, stc.), nonetheless, in the Pluribus design,

considerable attention has been given to maintaining as much isolation as practical so that one failure does not induce enother.

With the goels of flexibility and reliability in mind and with the price and size of minisomputers dropping, it was decided that the Piuribus should be built signs the lines of a minisomputers multiprocessor, or more generally, a multipresource (processors, memories, I/O shannels, etc.) system.

In considering which minicomputer might be most easily adaptable multi-resource structure, the to sommunication between the processor and its memory was of primary Several years ago maghines were introduced which Concern. sombined memory and I/O busses into a single bus. As part of this step, registers within the deviges (pointers, status, and control registers, and the like) were made to look like memory dells so that they and the memory gould be referenced in a homogeneous This structure forms a very eleen and attractive architecture in which any unit can bid to become master of the bus in order to communicate with any other desired unit. One of the important feetures of this structure is that it made memory acquesing "public"; the interfece to the memory had to become asynchronous, cleanly (solable electrically and mechanically, and well documented and stable. A Characteristic of this architecture is that all references between users are time-multiplexed onto a single bus. Conflicts for bus usage therefore establish an ultimate upper bound on overall performance, and attempts to speed up the bus eventually run into serious problems in arbitration.

In 1972 a new samputer **WAB** introducedesthe Lookheed SUEpawhish follows the single bus philosophy but derries it an important step further by removing the bus erbitration logic to a module separate from the processor. This step permits one to consider configurations embodying multiple processors, as well se multiple memories and I/O, on a single bus. It also permits busses which do not include any processor at all. The processor used in the SUE computer is a compact, relatively inexpensive (approximately \$600 in quantity), quite alow processor microcoded inner structure. Table 2=5 shows some characterizeics. Its slowness and ghosphess, of course, 90 together and since in a modular multi-processor, increased bandwidth is achieved merely by adding more processors, the weak/cheep progessor has the advantage of allowing smaller steps to be taken along the cost/performence curve.

16+bit word

8 General Registers

3.7 microseconds edd or load time

Misrosoded

Two words/instruction typical

8=1/2" x 19" x 18" shasa(s

64K bytes addressable by a

200 ne minimum bue avele time

850 ns memory syste time

425 he memory agrees time

Table 2-5 Sue Computer Cherectaristics

Several components of the SUE computer were adopted for the Pluribus system, in particular the bus, the processor, and the bus arbitration logic.

2.3.5.1 Hardware Strugture of the Plurfbus IMP

Reliability was a main consent in planning the hardware architecture. The goal was to provide hardware which sould be exploited by the program to survive the failure of any individual component.

The hardware consists of asynchronously and independently functioning communication busees, soupled together, physical point of view, the SUE chassis represents the basis construction units it incorporates a printed circuit back plane which forms the bus into which 24 eards may be plugged. From a incleal point of view this cheesis includes a bus which provides a common connection among all units olugged into the chassis. All specially designed pards as well as all Lockheed-provided modules plug into these bus chassis. With this hardware, the terms "bus" and "chessis" are used somewhat interchangeably, but we will commonly sell this standard building unit a "bug." Each bus requires one card which performs erbitration. A bus can be logically extended (via a bus extender unit) to a second bus if additional pard spage is required; in such a case, a single bus arbiter controls access to the entire extended bus,

One can build a small multiprocessor just by plugging several processors and memories (and I/O) into a single bus. For larger systems one quickly exceeds the bandwidth capability of a single

bus and is forsed to multimbus architecture. Please refer to figure 2-13 for the following discussions.

The functional units of the system (processors, memories, I/O controllers, and special devices) are distributed on these busses in such a way that units which must communicate frequently are placed on the same bus, whereas units which communicate less frequently will in general be on different busses. Units on the same bus can thus communicate at high speed without disturbing the remainder of the system. When a unit on one bus must sommunicate with a unit on another bus, some interference occurs while both busses are momentarily involved in the interestion. Each bus, together with its own power supply and socience is mounted in its own modular unit, permitting flexible veriation in the size and structure of systems. There are processor busies each of which contains two processors, each in turn with its own local 4K memory which stores frequently run and recovery-related code. There are memory busing to house the segments of a jarge memory common to ell the processors. Finally, there are I/O busses which gontrollers as well as certain central resources such as system clocks and special (priority-ordered) task disbursers which replace the traditional priority interrupt system.

The design is highly modular and permits systems of widely varying size and performance. In the interests of clarity,

Section 2 Design and Implementation

Figure 2-13; Plumibus Architecture

however, we will describe that design in terms of a particular system containing 14 processors==a system designed to have ten times the trafficehandling capability of the 316 IMP.

DRAFT

Referres. A central notion in a parallel system is the idea of a "resource", which we define to mean a part of the system needed by more than one of the parallel users and therefore a possible source of dentention. The three basis herdware resources are the memories, the I/O, and the progessors. It is useful to donsider the memories, furthermore, as a collection of resources of quite different character; a program, queues and variables of global nature, local variables, and large areas of buffer storage. The multiprocessor system is therefore in reality a multi-resource system, as mentioned above. The basic idea of the system is to provide multiple copies of the vital resources in the hope that the algorithm can run faster by using them in parallel and can survive fellures of apples of any of the resources. The number of copies of the resource which are required to allow sansurrent apprecian is determined by the speed of the resource and the frequency with which it is used. An additional edvantage of multiple sopies is reliability; if a system contains a few apere copies of all resources, it can continue to operate when one sopy breeks.

It may seem populier to think of a processor as a resource, in fact in a Piuribus IMP system the parallel parts of the algorithm compete with each other for a processor on which to run. Indeed, a novel feature of the Pluribus IMP design is the consistent treatment of all processors as equal units, both in the hardware and in the software. There is no specialization of processors for particular system functions, and no assignment of priority among the processors, such as designating one as meater. Net only the IMP application lab but also the multisprocessor control and reliability jobs are distributed among the processors so that all jobs are uniformly treated. The processors are viewed as a resource used to advence the algorithms the identity of the processor performing a particular task is of no importance. Programs are written as for a single processor except algorithm includes interlocks necessary to insure multiprocessor sequentiality when required. The software thus consists of a single conventional program run by all processors. Each processor has its own local gopy of about one quarter of this program and the remaining three querters is in commonly acquesible memory,

<u>Processer'husses.</u> The bus used in the Fluribus can up to four processors. As more are added, contention for the bus increases, and the performance increment per processor drops. The Pluribus IMP uses two processors per bus, which loses elmost nothing in processor performance. When a processor makes access to shared memory via the switching arrangement, that access inquestions due to contention and delays introduced by the intervening switch. In the IMP application, some parts of the program are run very frequently and other parts are run far less frequently. This allows a significant advantage to be gained by the use of private memory. A 4K local memory containing an individual sopy of the frequently run code is associated with each processor on its bus. This allows feater access to this "hot" code; the local memories all contain the same code. In the IMP application, the ratio of accesses to local versus shared memory is better than three to one.

Shared memory conteins program, message buffers, clobal variables, etc. Suffer requirements, of source, vary depending on the site. In a requier IMP, 40K words of common memory are necessary. Four memory units operating concurrently are required in order to hold processor contention to an acceptable level in the 14-prosessor system. Since the bus is considerably faster than the memories, two lociosi memory units may be placed on each shared memory bus with almost no interference. Two such memory busses are thus required.

I/Q.buggas. The I/O system consists of more standard busses, In the case of the 14-processor system, two such busses silow for the necessary bandwidth and provide redundancy for reliability

purposes, Into these busses ere plugged sands for each of the verious types of I/O interfeces that are required, including interfeces for modems, terminels, host computers, etc., as well as interfeces for stenderd peripherals. The I/O bus also houses a number of special units including (1) a dioak which serves roughly the same functions as the clock in the 516/316 IMP; (2) a shecksum/block transfer eard which flows a block of memory through itself computing a checksum as it does (used to checksum critical each, to compute end-to-endechecksums, etc.); (3) a special handwere task disbursing unit known as a Pseudo-Interrupt Device (PID) disquesed further below; and (4) a "reload" gard which monitors up to eight communication lines, watching for specially formatted reload massages from the network and processing them should eny errive.

Interesting everage. To adhere to the requirement that all progressors must be equal and able to perform any system task, busses must be connected so that all processors gan access all shared memory, so that I/O can be fed to and from shared memory, and to that any of the progressors may control the operation and so that any of the progressors may control the operation and some the status of any I/O unit.

A distributed inter-gammunization acheme was showen in the interest of expandebility, reliability, and design simplicity. The kernel of this scheme is called a Sus Coupler, and consists of

two gards and an interconnecting cable. In making connections between processors and shared memory, one gard plugs (nto a shared memory bus, the other into 4 processor bus, Similar connections are made for every processor but to every shared memory bus. the processor requests a cycle within the address range which the Sua Coupler recognizes, a request is sent down the cable to the memory and, which then starts contending for the shared memory bus. When selected, it requests the desired dysle of the shared memory. The memory returns the desired information to the Bus Coupler, which then provides it to the requesting processor, which, except for an additional delay, does not know that the memory was not on its own bus. The Bus Coupler also does address Since a processor can address only 64K bytes (16-bit address), and lince we wished to permit. multiprocessor configurations with up to 1824X bytes (28ebit address) of shared memory, a mechanism for address expansion is required. The Bus Coupler provides four independent Skebyte windows into shared memory. The processor can load registers in the Bus Coupler which provide the highworder bits of the shared memory address for sach of the four windows.

Given a Bus Coupler connecting each prognagor bus to each shared-memory bus, all processors can access all shared memory.

I/O devices which do direct memory transfers must also access these shared memories. These I/O devices are plugged into as many

I/O busses as are required to hendle the bandwidth involved, and bus couplers then gonnect each I/O bus to each memory bus. Similarly, I/O devices also need to respond to processor requests for sation or information; in this regard, the I/O devices act like memories and Bus Couplers are again used to gonnect each processor bus to each I/O bus. The path between processor busses and I/O busses is also used to allow processors to examine and gontrol other processors for startup and trouble situations.

2.3.5.2 Software Structure of the Pluribus IMP

The problem of building a packeteswitching storeeendeforwerd demmunications processor lends itself especially well to parallel solution since packets of date can be treated independently of one enother. Other functions of the IMP program such as general housekeeping, reutine computations, reliability tasks, etc., den also be easily performed in parallel. The structure chosen works as follows: first, the program is divided into small pieces, called strips, each of which handles a particular espect of the lob. For example, one strip handles special routing massages from neighboring IMPs, another handles input from a logal host, and others handle further I/O and housekeeping functions. When particular task needs to be performed, for instance upon reseipt of a massage over a communications circuit, the name (number) of the appropriate strip is put on a queue of tasks to be run. Each

processor, when it is not running a strip, repeatedly sheeks this queue, when a strip number appears on the queue, the next available processor will take it off the queue and execute the corresponding strip. The program is broken into strips in such a way that a minimum of context saying is necessary.

Strips have different levels of importance. Dete coming in over a high-speed communication girquit must be serviced more repidly then date coming in over a Teletype-speed line. The number essioned to each strip reflects the priority of the task it performs, when a processor checks the task assignment queue, it takes the highest priority job then everlable. Since ell processors assess this queue frequently, the contention for it is very high. For efficiency, therefore, a special herdwere devise, the Pseudo Interrupt Device, was designed to serve as a task queue. A single instruction allows the highest priority task to be fetched and removed from the queue, Another instruction allows a new task to be put onto the queue, All contention is arbitrated by standard bus logic hardware.

The length of strips is governed by how long priority tasks can wait if all the processors are busy. The worst case arises when all processors have just begun the longest etrip. In the IMP application, the most urgent tasks can efford to wait a maximum of 486 microseconds. Therefore, strips must not be longer than that,

(Of course, a strip might be longer if it is run infraquently and if the urgent tasks do not have absolute time requirements. That is, one might build a statistically acceptable set of strip lengths.)

An inherent part of multiprocessor operation is the locking of critical resources. This is the mechanism by which the algorithm enforces sequentiality when it is needed. Our system uses an uninterruptable loadeendedlear operation (load an accumulator with the contents of a memory location and then often the location) as its primitive locking facility (i.e., as the necessary multiprocessor look equivalent to Dijkstra semephores). To avoid deadlocks, we assign a priority ordering to our resources and arrange that the software not look one resource when it has already locked another of lower or equal priority.

Uning these fee; ities and techniques; the logical atructure of the IMP progrem is similar to that described for the 516/316 machines. In addition there is a substantial new progrem segment devoted to maintaining the hardware and keeping the program running in the face of hardware failures.

2.3.5.3 PLURIOUS IMP Reliability

Computer reliability is a common, serious, and difficult problem which has been approached in many ways. For critical

applications (*,9,, space exploration), large amounts of money are spent to overcome such apparently trivial weaknesses approblematical power supplies and connectors. Although a greet deal of attention is given to tailoring samputers to particular job environments, the commercial world of computer menufacturers has provided no adequate answer to the reliability problem.

The notions of faultetolerant and faileastt systems have been around for a number of years and because the reliability is such a crucial issue in a communications network, it was decided that some of these ideas should be exploited in the dealgn of the PLURIBUS IMP.

The reliability goal of the PLURISUS IMP is not that the system should never break, but rather that it should requerate automatically within seconds or minutes from most troubles and that the probability of total failure should be dramatically reduced over traditional machines, say to once a year or less. The system should survive not only transient failures but also solid failures of any single component. It is not necessary to operate correctly most of the time so long as outsides are infrequent, kept brief, and fixed without human intervention. Outsides of a few seconds are tolerated easily, and outsides of many seconds, while deusing the particular node to become temporarily unusable, will not in general Jeopardize operation of the network as a whole.

hapter III Section 2 Design and Implementation

Achieving this sort of reliability requires hardware that will survive any single failure, even a solid one, in such a way as to leave e potentially runnable mechine intact (potentially in that it may need resetting, reloading, etc.), Second, it requires all of the facilities necessary to survive any and all transfents stemming from the feiture and to edept to running in the new hardware configuration. To provide adequate hardware, extra sopies of every vital hardware resource are included. Sufficient isolation is provided between the copies so that any single component failure will impair only one copy.

Approximate bardware. It is not sufficient merely to provide duplicate copies of a pertiquiar resource; it is necessary to assure that the sopies are not dependent on any sommon resource; Thus, for example, in addition to providing multiple memories, there are multiple busses on which the memories are distributed; Furthermore, each bus is not only logically independent, but also physically modular. The chaesis, with its own power supply and cooling, is built into an integral unit which may be powered down, disconnected, and removed from the reak for servicing or replacement while the rest of the machine continues to run.

All central system resources, such as the real time clock and the PID, are duplicated on at least two separate I/O busses. All connections between bus pairs are provided by separate bus

souplers so that a doupler failure can disable at most the two busses it is connecting; all other interconnections between busses are unaffected.

When a particular communications circuit is deemed critical, it is connected to two identical interface units (on separate I/O busses), either of which may be selected for use by the program, when the extra reliability is not worth the extra cost, the line is only singly connected.

In order for the system to adapt to different hardware configurations, facilities have been provided which make it convenient for the seftwere to search for and locate those resources which are present and to determine the type and parameters of these which are found.

To allow for active fellures, all bus couplers have a program-controllable switch that inhibits transactions via that doubler. Thus, a "malicious" bus may be effectively "amputated" by turning off all souplers from that bus. These switches are protected from capricious use by requiring a password. Naturally an amputated processor has no access to these switches.

Software an experience any single component failure and still present a runneble system. One must assume that as a consequence

of a failure, the program may have been destroyed, the processors helted, and the hardware put in some hung state needing to be reset. Three broad strategies have guided the means used to restore the algorithm to operation after a failures keep (taken) almost, worry about redundancy, and use watchdog times throughout.

SIMPLICITY: First, all progessors are identical and equal; they are viewed only as resources used to advence the algorithm. Each is able to do any system tasks none is singled out (except momentarily) for a particular function. A consequence of this is that the full power of the machine can be brought to bear on the part of the algorithm which is busiest at a given time. A further consequence is that should any processor fell, the rest will continue to perform the mecassary tasks, albeit at reduced depactty.

A second system sharesteristic which gross from a desire to keep things simple is passivity. The terms active and passive describe communication between subsystems in which the receiver is expected to put saids what it is doing and respond. The quicker the required response, the more active the interaction. In general, the more passive the dommunication, the simpler the receiver can be, because it can weit until a donvenient time to process the communication. Neither the hardware interfaces nor

other processors tell a processor what to do; rather, tasks to be done are posted in the PID and processors ask the PID what should be done next.

There are some agets to such a passive system; first, the resulting slower responsiveness has necessitated additional buffering in some of the interfaces; second, the program must regularly break from tasks being executed to check the PID for more important tasks. The elternatives, however, are far worse, in a more eative system, for example one which uses classical priority interrupts, it is difficult to decide which processor to switch to the new tasks. The possibilities for deadlocks are frightening, and the general mechanism to resolve them cumbersome,

As a third example of simplicity, the entire system is broken into reliability subsystems which are parts of the overeil system that verify one another in an orderly fashion. The subsystems are cleanly bounded with well-defined interfages. They are self-contained in that each includes a self-test mechanism and reset depolitity. They are isolated in that all communication between subsystems takes place passively via data structures. Complete interlocking is provided at the boundary of every subsystem so that the subsystems can operate asynchronously with respect to one another.

The monitoring of one subsystem by enother is performed using timer modules, as discussed below. These timer modules guarantee that the self-stest mechanism of each subsystem operates, and this in turn guarantees that the entire subsystem is operating properly.

REDUNDANCY: Redundancy is simultaneously a blessing and a curse. It occurs in the hardware and the software, and in both control and data paths. We deliberately introduce redundancy to provide reliability and promote efficiency, and it fraquently occurs because it is a natural way to build things. On the other hand the more existence of redundancy implies a possible disagreement between the versions of the information. Such inconsistencies usually lead to erroneous behavior, and often persist for long periods.

There are several methods of dealing with redundancy. The first and best is to eliminate it, and always refer to a single copy of the information. When we shoose not to eliminate it, we shock the redundancy and explicitly detect and correct any indonsistancies. It does not really matter how this is done as the system is recovering from a failure anyway. What is important is to resolve the inconsistency and keep the algorithm moving. Sometimes it is to difficult to test for inconsistancy; then

There is a uniform structure for implementing a TIMERS: monitoring function between reliability subsystems based watchdog timers. Consider a subsystem which is being monitored. Such a subsystem is designed to sysle with a sharesteristic time constant, and a somplete self-sconsistency check is included within every dydie. Redular passage through this cycle is therefore sufficient indication of correct operation of the subsystem, excessive time does by without passage through the evele, it implies that the subsystem has had a failure from which it has not been able to recover by itself: The mechanism for monitoring the evole is a timer which is reset by every passage through the sysie. There are both hardware and seftware timers ranging from five microseconds to two minutes in duration in the IMP system. Another subsystem monitors this timer and takes correctlye action if the timer ever runs out. To avoid the necessity for subsystems be aware of one amother's internal atructure, each subsystem includes a reset mechanism which may be externally activated. Thus, corrective action consists merely of invoking this reset. The reset algorithm is assumed to work sithough a particular incornation in code may fail because it gets demaged. In such a dame another subsystem (the gode checksummer) will shortly repair the damage.

The entire system agnaists of a ghain of subsystems in which each subsystem monitors the next member of the chain. Lower

subsystems provide and certify some important environmental feature used by higher level systems. For example, a low level sode tester checksums all code (instuding itself), insures that all subsystems are requiving a share of the processor's attention, and guarantees that locks do not hang up. It thus quarantees the most basic features for all higher levels. These will, in turn, provide further environmental features, such as a list of working memory areas, I/O devices, etc., to still higher levels.

Before they can work together to run the mein eystem, a common environment must be established for all processors. The process of reaching an agreement about this environment is galled forming a consensus, and the group of agreeing processors is known as the Consensus, An example of a task requiring consensus is the identification of usable sommon memory and the essignment of functions (sode, variables, buffers, etc.) to particular pages.

The Consensus maintains and sounts down a timer for every processor in the system in order to detect unaccoperative or dead processors. This monitoring meshanism includes releading the failing processor's local memory and restarting it. Reliance on the Consensus is vulnerable to simultaneous transient failure of all processors. For many cases (as for example when all of the processors helt), a simple reset consisting of a sneesecond timer on the bus and a 60 Hz interrupt routine suffices.

for more detastrophic fellures the machine can be reset, reloaded, and restanted directly from the Network Control Center, which depends on the continual presence of human operators for successful operation. It is correspondingly powerful, resourceful, and erratic in its behavior.

2.4 Subnet Protocols

We begin this part of our report with a brist summary of four important protocols: IMP=IMP, Routing, Source IMP=Destination IMP, and Host=Host.

IMP-to-IMP transmission control in which each physical inter-IMP circuit is broken into a number of logical channels. Acknowledgements are returned piggybacked on normal network traffic in a set of acknowledgement bits, one bit per channel, contained in every packet, thus requiring less bandwidth than sending each acknowledgement in its own packet. Normally the number of logical channels is eight per physical line. A pair of IMPs connected by a high delay line such as a transoceanic satellite line may, however, be configured to use sixteen logical channels on that line in order to utilize the line more fully.

Each packet is essigned to an outgoing logical channel and carries the add/even bit for its channel (which is used to detect duplicate packet transmissions), its channel number, and the exknowledge bits - one for each channel in the reverse direction. The transmitting IMP continually system through its used channels (those with packets associated with them), transmitting the packets slong with the channel number and the associated add/even

DRAFT

bit. At the receiving IMP, if the add/even bit of the received pasket does not match the odd/even bit essociated with the appropriate receive channel, the packet is accepted and the respine odd/even bit is complemented; otherwise the packet is a duplicate and is disparded.

Every packet arriving over a line contains acknowledges for sil channels. The ack bits are set up at the distant IMP when it copies its receive odd/even bits into the positions reserved for the acknowledge bits in the control portion of every packet trensmitted. In the absence of other traifie, the acknowledges are returned in null packets in which only the acknowledge bits contain relevent information (i.e., the channel number odd/even bit are meaningless; null packets are not acknowledged), When an IMP regeives a packet, it compares (bit by bit) the ecknowledge bits egainst the transmit odd/even bits. For each match found, the corresponding channel is marked unused, the corresponding waiting pasket buffer is distanded, and the transmit odd/eyen bit is complemented.

In view of the large number of channels, and the delay that is encountered on long lines, some packets may have to wait an inordinately long time for transmission. A one-character packet should not have to welt for several thousand-bit packets to be transmitted, multiplying by 10 or more the effective delay seen by the source. Therefore, the following transmission ordering agheme has been instituted; priority packets which have never been transmitted are sent first; next sent are eny regular packets which have never been transmitted; finally, if there are no new packets to send, previously transmitted packets are periodically retransmitted even when there is a continuous stream of new traffic.

Each packet is individually routed from IMP to IMP through the network toward the destination. At each IMP along the way, the transmitting hardware generates initial and terminal framing characters and chacksum digits that are shipped with the packet and are used for error detection by the receiving hardware of the next IMP.

Errors in trensmission can effect a packet by destroying the framing and/or by modifying the data content. If the framing is disturbed in any way, the packet either will not be recognized or will be relected by the receiver. In addition, the check digits provide protection against errors that affect only the data. The sheek digits can detect all patterns of four or fewer errors occurring within a packet, and any single error burst of a length less than twenty-four bits. An overwhelming mejority of all other possible errors (all but about one in two to the twenty-fourth) is also detected. Thus, the mean time between undetected errors in the subnet should be on the order of years.

DRAFT Section 2 Design and Implementation

III=336

Routing Algorithm. The routing algorithm directs each packet to its destination elong a path for which the total estimated transit time is smallest. This path is not determined in advance, Instead, each IMP individually decides onto which of its output lines to transmit a packet addressed to another destination. This selection is made by a fast and simple table lookup procedure. For each possible destination, an entry in the table designates the appropriate next log. These entries reflect line or IMP trouble, treffic congestion, and current local subnet connectivity. This routing table is updated whenever necessary, as described below:

Each IMP estimates the delay it expects a packet to encounter in reaching every possible destination over each of its output lines. It selects the minimum delay estimate for each destination and periodically passes these estimates to its immediate neighbors, Each IMP then constructs its own routing table by combining its neighbors estimates with its own estimates of the delay to each neighbor. The estimated delay to each neighbor is based upon both quaue lengths and the recent performance of the connecting communication gircuit. For each destination, the table is then made to specify that selected output line for which the sum of the estimated delay to the neighbor plus the neighbor's delay to the destination is smallest.

The routing table is periodically and dynamically updated to adjust for changing conditions in the natwork. The system is adaptive to the ups and downs of lines, IMPs, and congestion; it does not require the IMP to know the topology of the network. In perticular, an IMP need not even know the identity of its immediate neighbors. Thus, the leased circuits could be reconfigured to a new topology without requiring any changes to the IMPs.

Source interpretination in presents in order for a source Host to communicate with a destination Host, both source and destination IMPs must establish a research of the connection for that Host pair. This simplex connection, consisting of a Transmit Message Block at the source, and a corresponding Receive Message Block at the source, and a corresponding Receive Message Block at the destination, is greated, and later removed, using a special protocol which detects duplicate or missing messages. The connection is dissilowed if the Host/Host agests control mechanism does not permit they host pair to communicate. A pair of Hosts may communicate with each other only if they are members of the same logical subnetwork or if one is allowed to communicate with Hosts in a subnetwork of which the other is a member.

To insure that messages arrive at a destination Host in proper order, an independent message number sequence is maintained for each connection. A message number is assigned to each message

at the sounde IMP and this message number is checked at the destination IMP. Out of an eightwhit message number space, both the sounds and destination keep a small window (qurrently eight) of valid message numbers, which allows several messages to be in flight simultaneously. Messages erriving at a destination IMP with message numbers outside of the surrent window or with message numbers already merked as requived are duplicates to be discerded, The message number consept serves two purposes; it orders the messages for delivery to the destination Host, and it provides for the detection of duplicate and missing messages. The message number is internal to the IMP subnetwork and invisible to the Hosts.

A sequence control system besed on a single source/destination connection, however, does not permit priority traffic to go sheed of other traffic. More generally, a Most may wish to specify a particular treatment for each message; thus, a separate connection is created for each "handling type", Currently, there are two possible handling types, require (for high bandwidth) and priority (for low delay).

Singe message numbers and reserved storage are so orition; in the system, very stringent and sareful procedures were developed to account for a lost message. The source IMP keeps track of all messages for which a RFNM has not yet been received, and the destination IMP keeps track of the replies it either has yet to send or has already gent. When the RFNM is not received for too long (presently about 30 seconds), the source IMP sends a control message (using the same message number) to the destination inquiring about the possibility of an incomplete transmission. Depending on the state of the reply table and message window at the destination, it will respond with either an indication that the message was not received on that it is out of range, or with a deprese duplicate reply (RFNM). The source IMP continues inquiring until it receives a response. This technique generally insures that the source and destination IMPs keep their message number sequences synchronized and that any allocated space will be released should a message become lost in the subnetwork because of a machine failure.

A connection is terminated either after a prolonged period of inactivity (presently 3 minutes), or a somewhat shorter period of inactivity coupled with the need for the Message Block by some other connection, or by the need to resynchronize a message number sequence that has been broken. The special termination protocol can be initiated by either the source or the destination in the first two cases above, or by the source in the third case, upon the receipt of an Mout of range? response to an incomplete query, upon closing a connection, both source and destination release all resources held or allocated for thet connection.

There is a facility outside of the normal Host/Host connection mechanism for sending and receiving a stream of "raw packets". These messages are identified by a special Host=IMP and IMP=Host code and bypass the connection mechanism. They are routed normally through the subnetwork, but no sequencing, error control, ressembly, or storage allocation is performed. Thus, they may errive out of order at the destination Host, some packets may be missing or duplicated, or packets may be thrown away by the subnetwork if insufficient resources are evaluable to handle them. No RFNMs or other messages are sent back to the source Host, about such rew packets.

Hasimbert. Restagais. A network working group, comprised of system oriented program designers from the research installations, was formed to define a protocol and develop techniques for communication between computers. Some constraints were imposed by the Host-IMP protocol and by the wide variance in specifics of the operating systems involved. A Host-Host protocol was specified, detailing procedures for establishing "logical" connections between a pair of Host computers wishing to communicate. The Host-Host protocols are implemented as a network control program (NCP).

The NCP is typically implemented as part of the executive system of the Host computer. The NCP is responsible only for

communicating with the IMP and with the NCP in other Hosts. Its primery functions are establishing, breaking, and switching connections between user-level progresses in a pair of Hosts, and controlling the flow of byte-streams transmitted over the established connection. System primitives are provided at the user-level interfece whereby user-level progresses direct the NCP to establish an terminate connections with user-level processes in other Hosts. Thus user-level progresses maintain control over the timing of transmission sequences and over the interpretation of the contents of the byte-streams transmitted.

The communications subnetwork and the Hostamost protocol provide a framework or set of building blocks for the development of user-level sepablisties. Several function-oriented protocols for the more common uses of the network have been specified by the network working group. Each protocol is implemented as a pair of cooperating user-level processes; a "user" process and a "server" process. The "user" process appears as a subsystem to the network user. When setivated by a user at a local Host, the "user" process interacts with the NGP to establish a sennection with the derresponding "server" process in a remote Host. When using these subsystems, the functions perfermed by the sammunications subnet and the NGP are transporent to the user.

A user can loo into a local host and activate a subsystem called TELNET which performs the functions necessary to connect him to any other host on the network. Alternatively, a user can connect to a TIP to access any network feetility. For all practical purposes, a user will appear both to himself and to the remote host as a user directly connected to the remote host. Once a connection is established, he follows the conventions of the remote host for localing in, executive functions, and for use of subsystems.

The TELNET protocol specifies a "virtual network terminal" which resolves differences in echoing conventions, interrupt signals, device rates, character sets, etc. Thus almost any type of local terminal can be used with any remote host. Many TELNET subsystems perform other peripheral functions such as saving a transcript of a session in a local file, sending a command file in place of user-typed input, and reporting whether various hosts are or have been up.

Many ad hos protogols have been implemented to provide capabilities for transferring files between hosts, for remote job entry to batch systems, and for interactively using graphics systems. These implementations generally require explicit user intervention and have not solved the difficult problem of dealing with incompatible data streems. There is an onegoing effort

within the network working group to define network standard protocols for these jungtions which will handle the incompetabilities in a uniform manner. The protocols are in verying stades of development, including some experimental implementations.

It should be noted that the primitives previded by the NGP are evaluable to any user-level process. Thus any user can build his own network functions by explicitly interacting with the NGP primitives making them an integral part of his application.

2.4.1 IMPEIMP Protocols

In this section we discuss some of the issues in designing node-to-node transmission procedures, that is, the packet processing algorithms. We tough on these points only briefly since many of them are simple or have been discussed previously, Note that many of these issues occur again in the disquasion of source-to-destination transmission procedures.

2.4.1.1 Basta Concepts

Butlering..end..Ripsiteing. As we noted in disquesting memory requirements, the amount of node-to-node peaket buffering needs to equal the product of the disquit rate times the expected adknowledgment delay in order to get full line utilization. It may also be efficient to provide a small amount of additional buffering to deal with statistical fluctuations in the arrival rates, (.a., to provide queueing. These requirements imply that the nodes must do bookkeeping about multiple packets, which rates the several issues discussed next.

Error. Costrol. We have disquised many of the aspects of node-to-node error control above; the need for a pecket checksum, its size, the basis of the adknowledgment/petransmission system, the decision on whether the line is usable, and so on. These procedures are critical for network reliability, and they should

therefore run smoothly in the face of any kind of mode or directly failure. Where possible, the procedures should be salf-synchronizing; at least they should be free from deadlook and easy to resynchronize [24].

Storage Alleration and Flow Comtrol. Storage allocation den be fairly simple for the packet processing algorithms. The sender must hold a dopy of the packet until it receives an acknowledgment; the receiver can accept the packet if it is without error and there is an available buffer. The receiver should not use the last free buffer in memory, since that would cut off the flow of control information such as routing and acknowledgments. In accepting too many packets, there is also the change of a storage-based deadlock in which two nodes are trying to send to each other and have no more room to accept packets. This is explained fully in [19].

The above implies that the flow dontrol procedures can also be fairly simple. The need to buffer a direct can be expressed as a quantitative limit of a certain number of packets. Therefore, the node can apply a sus-off test per line as its flow control throttle. More stringent rules can be used, but may be unnecessary.

Prierity. The issue of priority in packet processing is quite important for network performance. First of all, the

concept of two or more priority levels for packets is useful in decreasing queueing delay for important traffic. Beyond this, however, cereful attention must be paid to other kinds of transmissions. Routing messages should go with the highest priority, followed by acknowledgments (which can also be piggy-backed in packets). Packet retransmissions must be sent with the next highest priority, higher than that for first transmission of packets. If this priority is not observed, retransmissions can be looked out indefinitely. The question of preemptive priority (i.e., stopping a packet in midetransmission to start a higher priority one) is one of a direct tradeoff of bandwidth against delay since circuit bandwidth is wested by each preemption.

Refer is is. There has been much thought given in the packet witching community to the proper size for packets. Large packets have a lower probability of successful transmission over an encor-prone telephone line (and this drives the packet size down), while overhead considerations (longer packets have a lower percentage overhead) drive packet size up. The delayslowering effects of pipelining become more pronounced as packet size decreases, generally improving storesendsforward delay characteristics; further, decreasing packet size reduces the delay that priority packets are because they are waiting behind full length packets. However, as the packet size goes down, potential

effective throughput also goes down due to overhead. Metgalfe has previously sommented on some of these points [28].

has been suggested that the ARPA Network peaket size is suboptime; and should perhaps be reduced from about 1888 bits to This is based on optimization of node buffer 250 bits. utilization for the observed traffic mix in the metwork. However, the relative cost of node buffer storage vs. circuits is possibly such that one should not try to optimize node buffer storage. The true tradecif which governs packet size might wall be efficient use of phone line bandwidth (driving peaket size larger) ve. delay characteristics (driving packet size smaller). If buffer storage is limiting, one should just buy more (up to the limits of the address space, of gourse). Further, it is probably true that ld one is erving for high bendwidth utilization, buffer size must be large. That is, high bendwidth utilization probably implies the use of large packets, which implies full buffers, when idle, the buffer size does not matter.

As noted above, the choice of packet size is influenced by many factors. Since some of the factors are inherently in conflict, an optimum is difficult to define, much less find. The current ARPA Network packet size of about 1000 bits is a good compromise. Other packet sizes (e.g., the 2000 bits used in several other networks) may also be acceptable compromises.

However, note that a 2000-bit packet size generally means a factor of two increase in delay over a 1000-bit packet size, because even high priority short packets will be delayed behind normal long packets which are in transmission at each node. The use of preemptive priority might make longer packet sizes efficient.

Other authors have been cited as recommending a minimum length "pecket" of about 2000 bits because they have concluded that most of the messages currently exchanged within banks and airlines fit nicely in one packet of this size. To clarify this point, we note that they use the term "packet" for the unit of information we call a "message" and thus are not actually addressing the issue of packet size. We discuss message size below.

2,4,1,3 Subsequent Modifications

We now take a close look at the algorithm used in the ARPA Network for IMP-to-IMP transmission gontrol. As has been noted elsewhere, the inter-IMP modem interface hardware has the capability of concreting sheeksums for outgoing packets and checking the checksums on incoming packets. This allows packets which are demaged in transmission to be detected and discarded without acknowledgment. Peckets correctly received are acknowledged. A good IMP-to-IMP transmission sontrol algorithm must detect errors, ecknowledge good transmissions, and provide

retrensmission in the event of errors. In addition, the IMP-to-IMP transmission control elgorithm is improved if it detects duplicates that are sometimes generated by retransmission. An elgorithm which performs all four of these tasks is described below (it is like the HDLC data communications standard with a label sequence number per logical channel and up to 8 logical channels).

A number of logical "channels" are maintained between each pair of IMPs. Consider only one channel to begin with, and further consider packet transmissions in only one direction on this channel. Of sourse, acknowledgments go the other direction en the channel. At both the transmit and regeive and of this channel a one bit sequence number (a kept. We gall this bit an odd/even bit. Both transmit and receive odd/even bits are initialized to be zero. Also, at the transmit and, a used/unused bit is kept for the ghannel. It is of gourse initialized to zero, meaning unused. When it is time to transmit a packet, a check is first made for the channels being unused. If it was previously unused, it is marked as used and the packet is transmitted. The state of the transmit add/even bit is included with the packet, When the packet arrives at the receiver, assuming the packet is regeived correctly, the packet's odd/even bit is checked against the receive odd/even bit. If they metch, the packet is accepted and the receive odd/even bit is complemented. Otherwise, the

packet would be ignored. In any game the receive odd/even bit is returned as an adknowledgment. At the transmitter, if acknowledgment bit does not match the transmit odd/even bit, pasket has been suscessfully sent and asknowledged and the packet gen be discarded; the channel marked unused, and the transmit add/eyen bit complemented. Otherwise the adknowledgment is a duplicate and is ignored. Suppose now a second copy of the packet arrives at the redelver, a packet which was sent before the first agknowledgment had a change to get back to the transmitter. When this packet arrives at the receiver, its odd/even bit does not metch the receive add/eyen bit and so that pasket is discarded as a dupilcate. Nonetheless, an asknowledgment is sent for the packet using the present state of the require odd/even bit. the acknowledgment gets to the transmitter, it does match the transmit odd/even bis, so the asknowledgment is a duplicate and is ignared.

The anknowledgment bit is the state of the receive odd/even 12 18 bit after it is complemented rether than before Hence the need for the "not match" rule when the complemented. adknowledgment arrives at the trensmitter.

Because of the potentially long distances between IMPs, one channel (a not enough to keep the intermIMP lines fully loaded, Therefore, eight logical channels are supplied between each pair of IMPs (32 are supplied between Satelities IMPs). It is not necessary to maintain ordering of IMPsto=IMP transmissions singe peaket ordering is performed at the destination IMP, This means that the transmit channels can be filled in any convenient order, and at the receive side, packets can be forwarded anwards as soon as they are correctly received regardless of the channel over which they arrived.

To avoid requiring Reparete packets for acknowledgments, adknowledgment bits are "piggy=backed" in pagkets going the other way on the line. In fact, all might receive add/even bits are transmitted with every packet going the other way. In the absence of any traffic coing the other way on the line, a pecket carrying only the eight acknowledgments is sent. In either case, the acknowledgments get back to the transmitter as fest as possible. Therefore, the transmitter knows vary soon whether a packet requires retrensmission or not, allowing the use of a minimal t (meout before retranamission. "Piggy=backing" 411 acknowledgments into every packet going the other way program bandwidth, iina bandwidth, and buffar space over a system which sends individual asknowledgments in individual packets,

In view of the use of a number of channels and the delay encountered on long lines, some packets might have to welt an inordinately long time for transmission. Traffic that is

essentially interestive should not be subjected to waiting for several thousand-bit packets to be transmitted, multiplying by ten or more the effective delay seen by the source. Therefore, the IMPs maintain two sets of queues for each output line, and service ell priority transmissions before any regular transmissions. Preemptive priority is not employed. The following transmission ordering scheme has been instituted; priority packets which have never been transmitted are sent first; next sent are any regular packets which have never been transmitted; finally, if there are no new packets to send, previously transmitted packets are periodically retransmitted even when there is a sontinuous stream of new traffic. Routing has priority over all packets.

Figure Zeiß shows the formet of packets as they appear on the interwIMP directs. Packets are permitted to be variable length up to a maximum of about 1888 bits. The IMPFs modem interface transmission hardware adds framing characters to each packet as it is transmitted onto an interwIMP ejecuit. At the front of the packet two characters (DLE and STX) are added indicating the beginning of the packet. At the end of the packet two characters (DLE and ETX) are added indicating the characters (DLE and ETX) are added indicating the end of the packet. The checksum is appended after the end of packet characters. The IMPFs modem interface reception hardware has the capability of detecting the start and the end of a packet from these freming characters thus freeing the program from the burden of detecting

packet boundaries. Between peakets the transmission hardware automatically denorates synchronizing characters (SYN) which the reception hardware automatically dispards.

is no restriction on the content of a packet. Arbitrary sequences of bits may be transmitted without restriction. This transparency is achieved by a method known as Differentiating. If the date in the packet itself contains a DLE-character (the character which introduces the packet start and packet and sequences), that DiExcharacter in the data is doubled by the transmission hardware. At the redeiver, the hardware collepses double DLEs in the data back into one; so there is somplete transparency on the inter-IMP shannels. This is a very important point. All too many networks require transmission to be limited to characters from a particular character set. Besides preventing the network's users from sending applitmany messages, the designers of such networks are themselves prevented from such things as loading programs over the network,

In addition to showing the part of the packet formet done with hardware, Figure 2014 also shows the channel field, "plugy-back" acknowledgment field, etc., mentioned in the preceding section. The portions of the packet which carry the end-to-end control information (e.g., destination address, message sequence numbers, etc.) are also shown.

Section 2 Design and Implementation

Figure 2014; Packet Format

Notice that much of the intermIMP communication algorithm is performed with herdwere, Software is perticularly bad at generating powerful checksums, againing for packet boundaries, and so forth, especially when such chores must be done on a character-by-character or bit-by-bit basis. Therefore, they are done with herdwere in the ARPA Network.

Consider Figure 2=15. Two-way single-peaket treffig flows between A and As and also between B and Bs, and is constrained by network topology to use the circuit between IMPs C and D. Suppose all the buffer storage in IMP C can become filled with packets on the output queue to IMP D, and all the buffer storage in IMP D can become filled with packets on the output queue to IMP C. In this agas, both IMP C and IMP D would be engaged in a direct confrontation in which both IMPs must lose all incoming packets (and acknowledgmants) as neither machine has any buffer space with which to receive inputs. We call this a store-end-forward lockup.

It is straightforward to prevent such storemand-forward lookups, and this is done in the ARPANET implementation. The key technique used is to guarantee that sufficient buffers are reserved so that it is always possible to input and to output one packet over each circuit. Thus, packets (and adknowledgments) gan always be passed between neighboring machines (albeit at a trickle, paphaps). In particular, in the IMP system, one buffer

Section 2 Design and Implementation

Figure 2-15: StoremendeForward Lockup

1

is elways allocated for output on each line, guaranteeing thet cutput is always possible; and double buffering is provided for input on each line, which permits all input traffic to be examined by the program, so that adknowledgments can always be processed, which frees buffers. Additionally, an attempt is made to provide enough stere-endeforward buffers so that all lines may operate at full depective.

2.4.1.4 Constusions

We conclude this section with some comments: 1) negetive ecknowledgments could be useful in activating dormant buffers more quickly, but they add complexity and are not used in the ARPANET; 2) more complex forms of storesends forward lookup than that given in the example above are possible, and while most are protected against in the ARPANET implementation, at least one rare case is not protected against.

2.4.2 The Routing Algorithm

2.4.2.1 Basic Concepts

The fundamental step in designing a routing algorithm is the choice of the control regime to be used in the operation of the algorithm. Noneedeptive algorithms make no real attempt to adjust to changing network genditions; no routing information is exchanged by the nodes, and no observations or measurements are made at individual nodes. Centralized edeptive algorithms utilize a central authority which distates the nouting decisions to the individual nodes in response to network changes. Isolated adaptive algorithms operate independently with each node making explusive use of local data to adapt to changing conditions. Distributed adaptive algorithms utilize internode sooperation and the exchange of information to errive at routing decisions.

2.4.2.1.1 Noneadaptive Algorithms

Under this heading come such techniques as fixed routing, it is a fixed routing, and rendem routing (also known as flooding or selective flooding).

Simple fixed routing is too unreliable to be considered in precise for networks of more than trivial size and complexity. Any time a single line or node fails, some nedes become unable to communicate with other nodes. In fact, networks utilizing fixed

routing always assume manual updates (as necessary) to another fixed routing pattern. However, in practice this would mean that every routine network component failure becomes a catastrophe for operational personnel, every site apanding frantic hours manually reconstructing routing tables.

At their best, in the absence of network component failure, fixed routing algorithms are inefficient. While the routing tables can be fixed to be optimal for some traffic flow, fixed routing is inevitably inefficient to the extent that network traffic flows very from the optimal traffic flow. Unreliability and inefficiency are also characteristic of two elternative techniques to fixed routing which fall under the heading of non-adaptive algorithms; fixed routing with fixed elternate routes and rendom routing.

Noneedaptive algorithms ere all extremely simple and san therefore be implemented at low cost. They are thus possibly suitable for hardware implementation, for theoretical analysis, and for studying the effects of varying other network parameters and algorithms.

In conditation, we do not recommend non-adeptive routing for most networks because it is unreliable and inefficient. Despite these drawbacks, many networks have been proposed or begun with non-adeptive routing, generally because it is simpler to implement

end to understand. Perhaps this tendency will be reversed as more information about other routing techniques. in published and so network technology generally grows more sophisticated.

2.4.2.1.2 Centrelized Adeptive Algorithms

In a centralized adaptive algorithm, the nodes send the information needed to make a routing desision to a Routing Control Center (RCC) which distates its decision back to the nodes for estual use. The advantages disimed for a centralized elgorithm are; a) the routing computation is simpler to understand than a non-centralized algorithm, and the computation itself can follow one of deverel well known algorithms; b) the nodes are relieved of the burden and overhead of the nouting computation; g) more nearly optimal routing is possible because of the sophistication that is possible in a gentralized algorithms and d) routing can be evolded.

Unfortunately, the processor bandwidth utilization at the senter is likely to be very heavy. The elsesical elgorithms that a centralized approach might use generally run in time proportional to N gubed (where N is the number of nodes in the network), while their distributed counterparts can run (through perallel execution) in time proportional to N squared. While it may be a seving to remove computation from the nodes, it may not

be possible to perform a subic computation on a large network in real time on a single computer, no metter how powerful.

The dialm that more optimal routing is possible with a dentralized approach is not true in practice. To have optimal routing, the input information must be completely accurate and up-to-date. Of course, with any realistic centralized elgorithm, the input date will no longer be completely accurate when it arrives at the senter. Similarly, the output date we the routing decisions we will not go into effect at the nodes until some time after they have been determined at the center.

Distributed routing elgorithms, whether fixed rendom, fixed alternate, or adaptive, may contain temporary loops, that is, a packet may traverse a somplete circle while the algorithm adapts to network shange. Proponents of sentralized routing often argue that such loops can best be avoided by sentralization of the domputation. However, because of the time lags cited above, there may indeed be loops during the time of propagation of a routing update when some nodes have adapted the new routes and other nodes have not.

A dentrolized routing algorithm has several inherent weaknesses in the updating precedure, the first being unreliability. If the RCC should fell, or the node to which it is sonnected goes down, or the lines around that node fell, or a set

of lines and nodes in the network fell so as to partition the network into isolated components, then some or all of the nodes in the network are without any routing information. Of course, several steps can be taken to improve on the simple centralised policy. First, the RCC san have a backup computer, either doing another task until a RCC failure, or also on hot standby. This is not sufficient to meet the problem of network failures, only local outsides, but it is necessary if the RCC computer has any appreciable failure rate. Second, there can be multiple RCCs in different locations throughout the network, and egain the extra computers can be in passive or active standby. Here there is the problem of identifying which center is in control of which nodes, since the nodes must know to which center to send their routing input date.

A releted difficulty with sentralized algorithms lies in the fact that when a node or line fails in the network, the failed demponent may have been on the previously best path between the RCC and the nodes trying to report the failure. In this case, just at the time the RCC needs routes over which to requive and transmit routing information, no noutes are evaluable; the evaluability of new routes requires the very change the RCC is unsuccessfully attempting to distribute. Solutions which have been proposed to solve this "deadlook" are slow, complicated, awkward, and frequently rely on the temporary use of distributed algorithms.

Finally, centralized algorithms can place heavy and uneven demands on network line bandwidth as shown in Figure 2-16, near the RCC there is a concentration of routing information, going to from the RCC. This heavy line utilization near the center means that centralized elgorithms do not grow gracefully with size of the network and, indeed, this may place an upper limit on the size of the network,

2.4.2.1.3 Isolated Adaptive Algorithms

One of the primary characteristics of an isolated algorithm which attempts to adept to changing sonditions is that it takes on the character of a heuristic process; it must "leapn" and "forget" various facts about the network environment. White approach may have an intuitive appeal, it can be shown simply that heuristic routing procedures are unstable and are therefore not of interest for most prectical network applications. The fundamental problem with isolated adaptive algerithms is that they must rely on indirect information about network conditions, since dach node operates independently and without direct knowledge of or communication with the other nodes,

There are two basis approaches to be employed, separately or in tandem, to the process of learning and forgetting. These approaches can be termed positive feedback and negative feedback. One way to implement positive feedback was suggested by Baran as

Figure 2+16 Centralized Routing Algorithms

part of his hotepotato routing dostrine. Each node increments the handover number in a packet as it forwards the packet. Then the handover number is used in a "backwards learning" technique to estimate the transit time from the surrent node to the source of the packet. Clearly, this scheme has drawbacks because it lacks any direct way of adapting to changes. If no packets from a given source are routed through a node by the root of the network, the node has no information about which route to shoose in sending a massage to that source. In general, as part of a positive feedback loop, the routing election must periodically try routes other than the current best ones, single it has no direct way of knowing if better routes exist. Thus, there must always be some level of traffic traveling on any route that the nodes are to learn about, single it is only by feedback from traffic that they can learn.

The other helf of an edaptive isolated algorithm is the negative feedback sysle. One technique to use here is to penalize the choice of a given path when a pecket is detected to have returned over the same path without being delivered to its destination. The relation of this technique to the exploratory nature of positive feedback is evident.

An adeptive isolated algorithm, therefore, has this fundamental weakness; in the attempt to adept heuristically, it

must osofligte, trying first one path and then another, even under stable network conditions. This oscillation violetes one of the important goals of any routing algorithm, stability, and it leads to poor utilization of network resources and slow response to shanging conditions. Incorrect routing of the packets during oscillation ingreases delay and reduces effective throughput correspondingly. There is no solution to the meldeta oneillation in such algorithms. If the agaillation is demped to be slow, then the routing will not adapt quickly to improvements and will therefore dealers nodes unreschable when they are not, with the result that suboptimal paths will be used for extended periods. If the escillation is fast, then subspined paths will elso be used much of the time, since the network will be shronically full of traffic going the wrong way,

2.4.2.1.4 Distributed Adeptive Algorithms

In our experience, distributed adaptive algorithms have none of the inherent limitations of the above algorithms; e.g., not the inherent unreliability and inefficiency of noneadaptive eigorithms, nor the unreliability and size limitations of sentralized algorithms, nor the inherent inefficiency and instability of isolated algorithms. For example, the distributed adaptive routing algorithm in the ARPANET has operated for five years with little diffiguity and good performance. However,

distributed algorithms do have some predtical difficulties which must be overcome in order to obtain good performance.

Consider the following exemple of a distributed adeptive algorithm. Each node estimates the "distance" it expects a packet to have to traverse to reach seek possible destination over each of its output lines. Periodically, it selects the minimum distance astimate for each destination and passes these estimates to its immediate neighbors. Each node then constructs its own reuting table by combining its neighbors" estimates with its own estimates of distance to each meighbors. For each destination, the table is then made to specify that selected output line for which the sum of the estimated distance to the neighbor's distance estimate to the destination is smallest.

Such an algorithm can be made to measure distance in hops (i.e., lines which must be traversed), delay, or any et a number of other metrics including excess bandwidth and reliability (of course, for the letter two, one must maximize rather than minimize). The above algorithm is representative of a class of distributed adaptive algorithms which we consider briefly in the remainder of this section. For simplicity of discussion we will assume that distance is measured in hops.

2.4.2.1.5 Routing Processing Goals

In this section we outline some of the desirable characteristics for the processing of routing information. There are six divergent goals for the routing algorithm in its task of accepting the input data about the network and generating the required output.

Simplialty. Simplicity is the goal we list first, because it assumes increasing importance as further requirements are pissed on the routing eigerithm and the trend to complexity grows. There are two distinst kinds of simplicity that are of great value here, and in any computer algorithm. First, it is almost essential that the routing program running in a network be simple enough for a man to understand what is happening in a given situation. This is desirable not from the point of view of keeping the merely to permit him to follow the operation of the but algorithm and find problems and suggest improvements in its This may sound trivial, but in a very large network performance. even the most basic measurements and the most observations are difficult to undertake. Thus it is useful if the routing eigorithm (seelf does not present further somplexities to the man trying to interpret its behavior. In prectice, it should be possible to understand or predict the behavior of the routing algorithm without difficulty. For example, orderedependent or non-deterministic rules may be too obsqure to follow, particularly in a natuork environment,

The second kind of simplicity that is desirable is simplicity in the design and structure of the algorithm, so that it can be goded in a small, simple program. It is more likely that the program will be efficient and reliable if it can be expressed simply to the computer. It is always good software design to write simple programs, and this is especially true in a network environment, where the programs must run continuously in regities, and run on many computers.

Reliability. It is oritical that the routing algorithm be reliable in the face of node and line failures. Such failures must be expected, and, indeed, when they occur the successful operation of the routing deliquistion is most essential. Therefore, the mementary or prolonged melfunction of any component in the network should not interfere with the steady, accurate process of calculating the best routes for traffig in the network. When parts of the network are isolated from each other, and when they are reunited, the changeover should be managed smoothly. If a node fails to receive one or two routing messages, or erroneous data, the global reliability of the distributed computation should not be affected.

Strady. Etate laiuties. A basis requirement of any routing algorithm is that, given a static set of inputs, it should arrive et a steedy state solution which is accurate and stable. This is a very elementary qual. but one that should not be overlooked. The operation of the algorithm should not be so approximate that outputs oscillate under static input conditions. For instance, random or heuristic elegatithms are undesirable for reason. Also, orderedependent algorithms and techniques with many sease may not always arrive at the same solution to a special given set of routing input data with slight changes in the spacifications. Attention to this goal is particularly important in the early stages of designing an algorithm. Then, when test cases are being thought up, it is essential that the algorithm arrive at an accurate and stable solution for all the simple static tests that can be devised.

This means, for instance, that in thinking about algorithms that should work for a network with a sonnectivity that is assumed to change, one should test the elgorithm on the following kinds of networks:

eingle node loop networks ster networks tree networks union of a loop and a tree union of two loops series-peralish combinations

and so on. The algorithm should also be tested with various traffic requirements, such as:

no traffic common with a line in common

and so on. Solving the routing problem for a static attuation is much easier than for a dynamically changing one. It is also much easier to verify that the routing actually performs as desired.

Admitsion of the output data outlined above is not a one-time entoutation. As we have pointed out, the routing algorithm must run continuously, and its input data may change at any time. This leads us to the next requirement, that the algorithm be quickly adaptive to changes in network topology and traffic patterns. This point is of central importance, since routing for fixed inputs is a trivial problem by somparison. The computation to determine whether any paths exist between a pair of nodes must be

mensitive to changes in the configuration of the network. The pircuits and nodes in the network may fell, isolating some nodes from others, or there may be an elternate path remaining to sonnect them. That is, traffic should be routed around lines and nodes as long as any path exists to the desired destination. This may even require that a packet be returned over a path it has already traveled. For instance, sonsider a pagket going the short way around a circle from its source to its destination when a line in the path shead breaks, The should turn around and go back the other way around the direle. Similarly, direcits and nodes may be added to the network, and an adaptive routing algorithm has the advantage that these changes may happen smoothly and without any human intervention in the operation of the network. The algorithm must also be mensitive to changing traffig levels and patterns, since these will affect the computation of the paths of low delay and high eaparity.

adaptation to Carroe a Reightty of Routing. One important appear of the requirement that a routing algorithm be adaptive to changing network conditions is that it implies a priority structure to the task processing in the node computer. Routing must always have higher priority than packet processing because it may be essential to change the noutes being used to some new routes, especially at moments of high treffic load. If routing cannot be recomputed because of higher priority tasks, then the

network performance will degrade significantly as congestion sats in due to outsofedate routing information. In practice, this means that?

- is The input of routing messages must always be possible,
 - a. Storede must be reserved for the messages (a reserve of the dommon free list is necessary since the node dennot know sheed of input time whether the next input
 - will be a packet or a routing message).
 - b. The input process must be able to run often so that messages are not lost.
- 2. The output of routing messages must siweys be possible.
 - a. Storage must be reserved for the messages (this can be dedigated fixed table storage).
 - b. The output process must be able to run often enough to send routing messages as required, and they should have higher priority than any other transmission.
- 3. The routing update computation must always be possible.
 - a. Storage must be reserved for the routing update (this elso can be dedicated tables).
 - b. The updating process must be able to run periodically, either on a clock interrupt or called by the input and output processes which have been guaranteed to run frequently enough.

Adaptacion. to Change. 2. Speed: Given that the routing algorithm is adaptive to changes in the network, there are further desirable characteristics we can list. It should adapt as quickly as possible. Once the pattern of traffic has changed in the network, the old paths may become suboptimal, and perhaps even dounter-productive, interfering with other useful traffic. The modes should decide quickly and in a harmonious, sooperative manner how to use the resources of the network. When the routing algorithm must adapt, it should do so smoothly, without escillation, and without creating undue conflicts in other parts of the network.

Giobal. Obtimality. There are other coals in addition to the local choice of paths of low delay and high throughput. It is important that the routing algorithm meet gertain global nequirements as well. For instance, the algorithm should lead the nodes in the network to a global optimum in utilizing shared resources. It is not enough for an individual node to find a good path to another node, it must do so in the light of the needs of other nodes in the network. It is possible to think up routing algorithms which stabilize at several operating points in a given situation, and only one of them is the global best use of nescures. A good routing algorithm noutes high bandwidth traffic to eshieve the maximum global flow. Suppose, as in Figure 2x17, it is desired to send 50 kilobits/second of traffic from node 1 to

Figure 2=17 Global Optimality

mode 4 and at the same time it is desired to send 50 kilobits/second from node 6 to node 5. If the traffic from 6 to 5 as well as the traffic from 1 to 4 must pass over the link from 6 to 5, the global flow is 50 kilobits/second. If, however, either the traffic from 1 to 4 or the traffic from 6 to 5 were to go should the other way vis the link from 2 to 3, the global flow would be 100 kilobits/second which is the maximum global flow under the desired traffic inputs. Of course, it is glearly best for the 1 to 4 traffic to go vie the link from 2 to 3 since that elso minimizes the global delay.

Extract. This suggests that the algorithm should be fair to competition for shered resources. While fairness san sometimes be quentified, it is often a subjective judgment, but it is vital that the routing elgorithm operating at one mode does not prevent enother node from gaining a share of network services. Consider Figure 2=18, in which it is desired to send 1 unit of traffic from node 2 to node 2°, from 1 to 1°, from 2 to 2°, and so on. Clearly the maximal global flow is attained under this desired loading by stopping all traffic from 8 to 8° since any traffic from 8 to 8° neduces the global flow. But this maximum flow assignment is unfair. Let f be the flow from 8 to 8° and g be the flow from 1 to 1°, from 2 to 8° and g be the flow from 1 to 1°, from 1 to 8° and g be the flow from 1 to 1°.

Figure 2-18 Fairness

Section 2 Design and Implementation

TF = Total flow # f + N*g

Max TF [##0, g#1] # N

Min TF [fels dwa] # 1

fairness demands that 8 be able to get agms traffic to 81, even if it degreeses the global flow. It might be sonsidered fair to give all modes an equal share of the available bandwidth, or to allocate bandwidth to each node in proportion to demand or eyellable capacity. For instance, one assignment which is cortainly fairer than 100, out is:

#81/2, gm1/2; TF#(N+1)/2,

A different approach (a)

/#1/N; dm(N=1)/N; TF#N#1+(1/N)

which is quite close to the maximum flow for large N. definition of "fair" has some merits and drawbacks, but the important point is that the routing algorithm should be designed to adhere to some fairness dogtrine, or else some traffic will be locked out.

2.4.2.1.6 Routing Performance Measures

In this section, we will take up the question of how to evaluate the performance of a routing algorithm. We will use the four basic factors underlying the performance of the network ex. a whole in considering the routing algorithm,

Delay. As we indicated in the section above on network delays, the most important point about delay usually is that interactive traffic experience the minimum delay possible, although spmetimes uniformity of delay is important. general considerations such as giving such traffic higher priority then bulk traffic, what effect can routing have on delay performance? One way to enswer this question is to review the components of delay introduced above, and to note what action the routing algorithm can take with regard to each. Before giving this list, it should be noted that the actual values of the veriables in question here may vary enormously from one network to another, and within a given network. Therefore, some of the considerations will pertainly outweigh others,

- In Propagation delays. The routing algorithm may keep track of the speed-of-light delays in the network, which may vary significantly if there are some long lines, and sertainly if there are setablish. Although these delays are fixed and beyond the control of the algorithm, it can shoose paths with the least propagation delay.
- Transmission delays. The routing algerithm may also keep track of the bandwidth of the streuits in the network, to know the transmission delays that peckets will experience. Again, these delays are not likely to be dynamic, but the routing computation can use feet lines and avoid slow lines where possibles

Section 2 Design and Implementation

3. Nodel processing delays. Here delay is more likely to have a large dynamic range. Some nodes may have different traific loads than others by several orders of magnitude, and the nodes themselves may have different capacities, so the input quausing delays for prodessing service may very. If this component can be significant, the routing election should monitor the values for processing delays at the nodes,

4. Queueing delay. The case of queueing delays on einquits is similar. There will most sertainly be some lines which are loaded more heavily than others, and the queueing delays are inversely proportional to line bandwidth, so long queues on slow lines lead to very long delays. In both of these cases, not only can the routing algorithm evoid long delays by choosing alternate routes, but it is also the major determinant of processing and queueing delays. That is, the routing algorithm may be structured to senso the buildup of these delays and shenge the routes being used eccordingly, explicitly to reduce the load on individual nodes or lines.

5. Retrensmission delays. This case is somewhat similar as well. The routing eigorithm should edept to any sizable number of retransmissions on any circuit, since they effect many of the other variables. By reducing the effective bandwidth of the elrouit, retrensmissions ingreese the processing and quauting

delays experienced by all packets. In this the transmission delay for a given dirguit is not strictly constant, This is especially true for matellite links used in broadcast mede.

Thraugheut. The next topic we will examine is that of network throughput as related to routing. We have stated that the important goal here is that large data transfers, as opposed to interective traffic, gain high throughput. A liet similar to that for delay can be proposed; these are issues which may be important for the achievement of high throughput, depending on the kind of network donalderedi

Circuit bandwidth. The routing algorithm must exceptain the effective bendwidth of the circuits in the paths that it chooses for high throughput traffic. This includes such effects use of the line by other modes, the deterioration retrensmissions, and so on. It is given that in this regard routing program is capable of much more than a passive measurement role. Since it is routing which determines at each node how much traff(s to send over a particular line, it is possible to requiste the flow over high throughput paths if desired. It possible, in order to provide as high a utilization rate as possible, to enqure that all traific is long packets and messages, although sending all interactive traffig by other routes penalize that traffic in terms of delays,

- 2. Node bandwidth. The same comments apply to node bandwidth,
- 3. Buffering. We will assume that routing and flow control ere different operations and are unfalated. Thus, the routing program need not be concerned with the existence of buffering at any level in the network.
- 4. Multiple paths. We have already mentioned the necessity for load-splitting in certain natwork situations. It is likely to be much more relevant for higherhroughput applications than for low-delay traifie, but the technique of using many paths to a destination is of fundamental importance.
- Cost. We now turn to a gonalderation of the relationship between routing and network cost. There are 3 basic network resources which the routing algorithm utilizes, and it can take several actions to utilize them effectively:
- Line bandwidth. A primary consideration here is that the routing algorithm elect the paths with the fewest intermediate nodes, thus minimizing the total line bandwidth utilization. This is an important argument, and the basis for severel algorithms based on shortest path routing. We should note that this goal may be in direct conflict with the goals of low delay and high throughput, though there are often seems, particularly in networks with uniform node and line characteristics, when shortest path routing is an excellent policy.

- 2. Node bandwidth. A second cost factor, similar to the first, is node bandwidth. Again, the souting algorithm which chooses short paths over long ones is at an advantage here. In both cases, the routing program would do well to seek underutilised resources rather than concentrating traffic on a few paths.
- 3. Node storage. The final cost factor is node storage, a specific instance of the higher cost of concentrated traffic. The routing algorithm has the capability of keeping the network queues as short as possible if that is a specific objective, Avoiding the inefficient use of storage in overlong queues may also tend to defer problems of congestion, which also make the network efficiency suboptimal.

A different sost consideration is that of network connectivity. Here too, the nouting algorithm is an important factor in determining gest, though in a more indirect faction. The cost of connecting the network is related to how adaptive the routing procedures are in preftice. For instance, some networks have been proposed with fixed routing matrices giving two elternate routes to each node. In this kind of a network, it is essential to provide enough somnectivity so that nodes are selded declared unreachable by the routing algorithm (this may happen even though a network path exists between them). Of course, this neighbors the cost of the network. A similar problem holds for

routing algorithms which adept elowly, or only with human intervention, or only with some given probability of accuracy, and so on. When we discuss area routing later in Section 4, it will be clear that the problem of rapid and accurate determination of reschability is an important and difficult problem in networks with hundreds of modes or more.

Reliability. The last performance measure we will discuss to network reliability. Routing on have several different kinds of effects on reliability. In terms of network use, the important measure is the frection of messages which are undelivered due to failures in the communications subnetwork. The routing program has as its main function in the network the efficient and reliable delivery of messages to their destinations. Despite all kinds of component failures, from lines to nodes to Hosts, the routing eigorithm should continue to deliver messages properly or report that they are undeliverable because the destination is unreachable or not functioning.

The perallel requirement conderning the reliability of network connectivity has been exemined in the section above en network dost. The appropriate measure here is the fraction of the time that the routing algorithm is in error conserning the reschability of some node.

A different set of issues snises in the relationship between the nouting program (tasif and network reliability, Given the sentral role of the routing program in any network, it is particularly important that routing never break down altogether, the way most systems, softwere and hardware, eventually do. The reason is clears if the routing in the network is incorrect or nonfunctioning, the network is completely unusable. This means that a different measure of routing performance is the percentage of network unavailability that is due to routing algorithm failures. In summary, the routing program must also be considered as another module of network software, with some given level of reliability, which is more sensitive in terms of network reliability than most other modules, because of its global impacts.

2.4.2.1.7 Routing Cost Measures

There are five besid dosts involved in the sontinuous operation of a routing algorithm. All these factors set to reduce the effective capabilities of the modes and times in the network with regard to the progressing of data packets. In other words, these costs are various kinds of network system overhead. First, there is the CPU utilization at each node in the network meeded to perform the deligation of the new best routes to all destinations. Second, there is the delay at each node associated with the processing of the routing information, introducing delays

in the processing of packets. Third, there is the storage needed at each node as a data base for the routing calculation, and for the exchange of routing information with other nodes. Fourth, there is the line bandwidth used for the exchange of routing information between nodes. Finally, there is the line delay seused by the fact that sometimes a routing message is being transmitted at a time when a data packet is queued. We will now consider each of these costs in more detail.

Line Bandwidth. The first nost feator to consider is the fraction of the evaluable line bendwidth needed to exchange routing messages between nodes. Clearly, this bandwidth is a function of the size of the routing message and its fraquency as previously stated:

BWCr a Brafr

Depending on the eigenithm, one may shouse also to make the bendwidth used for routing a fixed number of bits per second, regardless of available line bendwidth, or one may wish to keep the bandwidth used for routing below some acceptable overhead fraction of the line bandwidth. Thus, on slow lines, routing would be sent less often or in an abbreviated form. One may also shoose different priority strategies for the transmission of routing. As we have pointed out, routing messages are very important, and should probably take pregadence over most other

traffic. However, it may be desirable to apacify that some diseases of transmission take priority over routing messages. In this way, the routing messages need not introduce added delays to special high-priority messages, though they still represent a reduction in the effective bandwidth of the dommunications elective.

Line Delay. Next we exemine the added delays on singuita caused by the transmission of routing messages. Consider a single line with one routing message of length 5r sent with frequency from second, and a very light date traffic load. Then the probability of a data packet having to wait is the ratio of the time duration in which routing is being sent to the nouting period. For a given singuity bendwidth SWC, this probability of line delay is:

BWC+/BWC & Br*F+/BWC.

Given that a packet must wait, the average wait is half the meximum, or

ALD a Br/(2*8WC).

The expected line delay is the product of the probability of line delay and the average line delay:

ARPANET Completion Report
Chapter III Section 2 Design and Implementation

ELD # (BWC#/BWC)#ALD.

That is,

ELO = (BreBrefr)/(2+8WC#8WC),

This means that the delay to packets due to routing messages increases as follows:

ifnearly with increasing routing frequency, quadratically with increasing routing message length, quadratically with degreesing line bandwidth.

The reason that the analysis presented above is valid only for light traffic loads is that a queueing phenomenon square total network delay to increase nonlinearly with load. Therefore, the added traffic due to the presence of routing massages has a correspondingly greater effect in terms of delay at high loads.

Madal Bandwidth. The calculation of the best routes to elidestinations represents a steady, periodic demand for CPU
processing time. One can view the computation as a certain
percentage of overhead in the CPU bandwidth available for message
processing. We have gailed this fagtor BMPr, the fraction of the
bandwidth of the processor used for routing. In general, it has
two components, one based on the time taken to process routing
messages on each line, and the other a simple periodic components

BWPP & NLN#PP#FP + Pp

where NLN is the number of lines per node, Pr is the progressor time for routing messages, Fr is the frequency of routing messages, and Pp is the fractional overhead of progressor time for periodic processing of routing. For instance, if NLNEZ, Fr85 messages/seg, and Pps1%, then BWPr83%. If BWPr becomes too large, the processor becomes much less gost-effective in its primary role as message processor. Therefore, a major sost consideration in evaluating a routing algorithm is the number of CPU cycles per second it requires in each node of the network.

Madal Delay. Along with the reduction in effective node; processing power games the effect of delays in the processing of packets while the routing computation is proceeding, given that it takes priority over data processing. Suppose one routing message is received from each of NLN adjacent nodes with frequency from escapes per second. Then, given the time Pr and the frequency time Pp above, there are three cases to consider for nodel delays

- 1. NLNaPrafe >> Pp message processing dominates
- 2. NUMERPARE 44 Pp. periodic processing dominates
- 3: NLN+Pr+Fr # Pp factors have equal magnitude

If we assume that NLN+Pr+Fr 4 : (#total progessor bendwidth), and that packet processing time is infinitesimal, we can analyze each of these cases:

1. The probability that a packet has to wait is

NLN+Fr+Pr

and if it must wait for only one massage, it experiences an average wait of

Pr/2.

This means that the expected delay (s

NLN+FP+PP+PP+/2

which is a lower bound on the actual value. It can be shown with a more detailed energies that this is the expected value for delay under the condition that 1/(Fr#P#) >> NLN, which means that there are many more time periods in which the processor is able to process routing inputs than there are adjacent nodes to send the routing. In practice, this is the only assumption that makes sense, since otherwise SWPr begomes close to unity, and there is no bendwidth available for data traffice.

- 2. This case can be analyzed as line delay was analyzed above, based on the length and frequency of the periodic processing. Fp.
- 3. This case can be enalyzed on the basis of the first two cases, summing the effects of each of the terms.

Nodil Sterage. The final cost factor we will discuss is the storage required at each node to maintain the routing information. This does may vary greetly among various algorithms, depending on how much information is needed in the routing computation, and also on the method of exchanging routing messages. The node must save the incoming routing information in some fashion, then the routing computation may generate other, new information, and the transmission of routing messages to the adjacent nodes may call for still more storage for data. In a small communications progressor, memory is dedicated to a fairly small program, and message buffers. Any storage used for the routing calculation must be taken from the message buffer pool, either once and for all at design time, or dynamically. Therefore, the storage requirements of the routing algorithm represent atill another overhead factor.

2,4,2,2 Original Implementation

In this section we describe the routing algorithm originally installed in the ARPANET, and examine it from the point of view of the performance and gost criteria outlined above.

The original ARPANET algorithm can be summerized as follows. This algorithm directs each packet to its destination along a path for which the total estimated transit time is smallest. This path is not determined in advance. Instead, each IMP individually

decides which line to use in transmitting a packet addressed to another destination. This selection is made by a simple table lookup procedure. For each possible destination, an entry in the routing table designates the appropriate next line in the path.

Each IMP maintains a natwork delay table which gives an estimate of the delay it expects a packet to enquiter in reaching avery possible destination over each of its output lines. This table and other tables mentioned below are shown in Figure 2=19 as kept by IMP 2, for example. Thus, the delay from IMP 2 to IMP 5 using line 3 is found to be 4 in the Natwork Delay Table. Periodically, every 2/3 of a second, the IMP selects the minimum delay to each destination and puts it in the minimum delay table. It show notes the line giving the minimum delay and xeeps the number of the line in a table for use in routing packets. Also every 2/3 of a second, the IMP passes its minimum delay table to each of its immediate neighbors, that is, it sends the minimum delay table delay table out each of its phone lines. Of course, before the minimum delay table is transmitted to the neighboring IMPs, the IMP sets the minimum delay table is transmitted to the neighboring IMPs, the

Single all of the neighbors of an IMP are also sending out their minimum delay table every 2/3 second, with their own entry set to zero, an IMP reactives a minimum delay table from each of its neighbors every 2/3 second. These tables are read in over the

Figure 2-19 ARPANET Routing Algorithm

rows of the delay table as they arrive. The row to be written over is the row dorresponding to the phone line that the arriving minimum delay table came in over. After all the neighbors? estimates have arrived, the IMP adds the delay saved by the IMP itself to the neighbors? estimates. This is done by adding the IMP delay table, i.e., the contribution of this IMP to the total delay to each destination, to each column of the delay table. Thus the IMP has an estimate of the total delay to each destination over the best path to that destination.

In parallel with this computation, the IMPs also compute and propagate shortest path information in a similar fashion. This information is used only in the determination of connectivity. An upper limit of the number of lines in the longest path in the network is used as the cutwoif for disconnected or nonexistent nodes.

Now let us sonsider the performance of this elgorithm. First of all, it explicitly determines the connectivity of the network, since all IMPs are continuously exchanging the length of the shortest path from each IMP to each other IMP. Information travels at roughly 2/3 of a second per line, so that shanges in tapology are resognized by the whole network in a matter of a few seconds. This figure is probably asceptable if one assumes that the network connectivity does not change too often. Second, the

Section 2 Design and Implementation

algorithm also explicitly desculates the path of least delay, However, here the approximations due to the low frequency of routing update mean that the delay for traffic at one instant is a function of the traffic of several seconds before, This doubt potentially lead to oscillations and poor line utilization. The ARPANET algorithm attempts to head off this class of problems by biasing delay heavily toward the shortest path. That is, delay is measured by the number of packets on an output quoue, plus a fixed ingrament, so that even an empty quoue represents additional delay.

The algorithm has several faults, some of which are relatively simple to dure, and others which are more fundamental in neture. The strong bias in the eigorithm towards the shortest path is basically a good idea, and leads to stable flows near optimum values. However, the bias makes the algorithm somewhat insensitive to changes in traffic patterns, so that global optimization of delay and throughout is not likely as network loading increases. A second fault is that the election only maintains one route per destination, updated every 2/3 second, This means that no load-splitting is possible, at least not on a short term besig. The elegation could be modified to use one of several routes to each destination, with weights assigned to each, Further, the algorithm might maintain additional data on the loading of the various paths, to facilitate more rapid adaptation to changes in traffic. This change, combined with the expansion to several routes, might also lead to smoother and more uniform edeptation.

The original ARPANET routing algorithm was quite a good design in many respects. Perhaps its strongest point is that it is simple. The IMP does not have to know the topology of the network, or even the identity of its neighbors. When IMPs and itnes go down, the algorithm functions as usual, and the new routing information propagates through the network by a process of exchanges between neighbors. Therefore, the algorithm accres quite well in reliability. Although there are no explicit controls to ensure fairness to competition, the algorithm does relatively well in this estendary as well.

Finally, the algorithm is not a costly one in terms of the measures discussed above. The program in the IMP picks the minimum daily and hop counts from the routing messages received from each neighbor, for all destinations. Thus, the calculation is proportional to the number of IMPs in the network, and the number of lines donnected to each IMP. The routing computation takes up about 5% of the CPU bandwidth of the IMPs. The delay and hop information is packed into a single idebit word, so that the routing message sent out seek line consists of 64 words, one for each IMP in the network, plus some header information. This

amounts to less then 2% of the bendwidth of a 50 Kbs line. At these low bendwidth retes, added node deleys and line deleys are not appreciable. In addition to the storage required for sending the routing message out each line (one copy of the message is shared by all lines), the IMP receives storage for receiving routing messages from each of its lines. These tables, together with its own directory of the best line to each destination, amount to about 3% of the core storage on an IMPs. In summery, the IMP routing algorithm is a simple, inexpensive algorithm which performs well in steady state, and in reacting to small changes in traffic.

2.4.2.3 Subsequent Modifications

One problem with distributed algorithms is that they are slow in adepting to some kinds of thends; in particular, the original ARPANET algorithm resuts quickly to good news, and slowly to bad news. If the number of hops to a given node degresses, the nodes soon all agree on the new, lower, number. If the hop count increases, the nodes will not believe the reports of higher counts while they still have neighbors with the old, lower values, Another point is that there is no way for a node to know sheed of time what the next-best or islimback path will be in the event of a feilure, or indeed if one exists. In feet, there must be some

a shange in the network and the edeptation of the routing algorithm to that shange. This time depends on the size and shape of the network.

In 1973, the decision was made that the routing algorithm should continue to use the best route to a given destination, both for updating and forwarding, for some time period after it gets worse. That is, the algorithm should report to the edjagent nodes the gurrent value of the previous best route and use it for routing packets for a given time (nterval. This technique was "held down". One purpose of held down is to distinguish between changes in the network topology and traffle negetaitate shanging the choice of the best route, and these changes which merely affect the characteristics of the route, like hop count, delay, and throughput. In the case when the identity of the path remains the same, the mechanism of hold down provides on instantaneous adeptation to the changes in the characteristics of the Paths certainly, this is optimal. When the identity of the path must ghande, the time to adapt is equal to the absolute minimum of one network response time, while the other nodes have a change to react to the worsening of the best path and to decide on the next best path.

The routing algorithm is extremely important to network reliability, since if it malignations the network is useless,

Further, a distributed routing algorithm has the property that all the nodes must be performing the routing computation correctly for the algorithm to be reliable. A local failure can have global consequences; e.g., one node announcing that it is the best path Routing messages between nodes must have shecksums to all mades. discarded if a checksum error is detected. All and must be routing programs must be checksummed before every execution to verify that the gode about to be run is correct. The checksum of the program should include the preliminary checksum somputation itself, the routing program, any constants referenced, and anything else which could affect its successful execution. time a checksum error is detected in a node, the node should immediately be atopped from participating in the routing computation until it is restored to correct operation again.

The new concept introduced in the ARPANET, known as "held down", facilitates the process of adepting quickly, uniformly, and without oscillation. This technique is of gentral significance to routing design, and many aspects of its use are examined and explained. These shanges make it work considerably faster and also more efficiently. The provious approach, based on path length, required routing messages to be exchanged which contained hop counts to all nodes. The need for this mechanism is eliminated through the use of hold down. Finally, later sections consider the application of this technique to the problems of traffic assignment.

The original ARPANET routing technique has a flaw that it shares with all algorithms built on a repeated distributed minimization or maximization. The algorithm operates in such a way that the hop count to a given node increases smoothly at greater distances from that node. No node even has a hop count to a given destination which differs by more than one from the minimum dount held by any of its neighbors. What this means in terms of adeptation is that

the reachability algorithm reacts very guickly to good news, and very slewly to bed news.

If the number of hops to a given node degreeses, the nodes soon all agree on the new, lower, number. If the hop dount increases, the nodes will not believe the reports of higher counts while they still have neighbors with the old, lower values. The simple case of a line of nodes shows how this happens. The four cases in figure 2×20 below show how the hop counts shange when the actual number of hops increases and decreases, for the best and worse sames of sequencing among the nodes. Clearly, the order in which the nodes compute and exchange routing is critical here, and we will assume that there is some definite order in which the nodes periodically perform the calculation and updating. We assume that each node in the sequence finishes the calculation, sands out the new tables, and the information errives at the nodes adjagent to that node, all before the next node in the sequence begins its computation.

DN MAX MAX MAX MAX DN MAX MAX MAX

UP 1 2 3 4 UP 1 MAX MAX MAX

1 2 MAX MAX

1 2 3 MAX

1 2 3 4

a. Left-to-right sequence, b. Right-to-left sequence, best case, node & comes up worst sage, node & comes up

UP 1 2 3 4 UP 1 2 3 4

DN 3 4 5 6 DN 3 2 3 4

5 6 7 8 5 4 3 4

7 6 5 4

9 8 7 6

B P DI 11 XAM XAM XAM XAM

...

MAX MAX MAX MAX

c. Left-to-right sequence, d. Right-to-left sequence, best case, node 8 goes down worst case, node 8 goes down

Figure 2=28 How the Reachability Algorithm Adapts

In the light of this analysis, we can remark on the choice of the hop count mechanism for the ARPANET reachability algorithm, is could have been built around minimum delay, which was being computed enyway for route selection, instead of around minimum hops. However, the value of delay corresponding to unreachable nodes is much higher than MAX, since hop counts increase by a single unit for each line, and delay has a large dynamic range (4 to 24 in the original ARPANET implementation, with all line depedities equal). Thus, using a high value of delay to detect unreachable nodes would take proportionately more time, given the linear rise time characteristic of the process. This time would have been so long, on the order of minutes, that a delay-based reachability elgorithm was not practical.

One can investigate the function served by rise and fell times (times required to adapt to or change) in a routing algorithm, in order to determine whether these delays are necessary, and if bounds can be placed on them. To begin with, it is easy to understand why a distributed algorithm has a zero fell time. If a node requires information about a better path, it will use it right away. As we shall see, however, this information may not be assurate or languithed, and therefore acting on the information may be a poor decision. On the other hand, it is also may to see why the distributed approach leads to a long rise time. When it discovers that the current best path has grown more

costly to use, it has no secondebest path readily svaliable. This is an instance of solving a global problem on the basis of local estimates. The node must choose the best path to a destination, and the only direct information it has is about the lines to which it is connected. The nouting information, which is about paths, is always slightly out of date and thus inacqurate. Further, the simple distributed routing algorithm only sends routing data about the best path, not all possible paths. Even if such information were aveilable to the routing decision process, it still requires the cooperation of the whole network, and thus some time and many neuting messages, to determine the identity of the nextwost path after the best path degrades. We next give some examples to illustrate this process.

Consider the networks shown in Figure 2=21. In all these networks, we are gongerned with a path from node 1 to node 3. Sugh a path exists, two hops long, in all three networks shown, when all components are up. In Figure 2=21s, a line has gone down, and no path exists from 1 to 3. The same is true in Figure 2=21b, where node 3 has gone down. In both these cases, the routing information held in the syste of nodes 1,2, and 4 will directly around, with the hop counts to node 3 slowly counting up. Initially node 2 is 1 hop away from node 3, and nodes 1 and 4 are 2 hops away. After the component fails, node 2 thinks there is another path through node 1 or node 4, and sets its hop counts

to 3. Likewise, node 1 counts to 4. Depending on the sequence in which the nodes send routing, the modes will begin to sount until the hop sount reaches MAX.

The other cases shown in Figure 2-21 provide a possible interpretation of this counting sequence. In the network of Figure 2-21s, the line from 2 to 3 is down, but an elternetive path of 3 hope exists. This is not true for the same network in Figure 2=21d, where node 3 is down. Similarly, in Figure 2=21e, a path of 4 hops exists and is the shortcet path, while in Figure 2=21f, there is no path to node 3. Obviously, there are an infinite number of such examples, each constructed upon the last, When a shartest path of H hops is disrupted, the routing eigerithm which is operating in the metwork at large must investigate for the existence of a new path, Conceptually, this can be done by lacking for a path of langth H, and then (if unsuggessful) for one of length Help H+2, H+3, ..., MAXw2, MAXw1, At that point the algorithm can conclude that no path exists to the given node. As indicated by the examples in Figure 2-21, this search must involve all the nodes and lines in the network. There is no way for a node to know sheed of time what the next sheet or fellaback path will be in the event of a failure, or indeed, if one exists. This line of reasoning indicates that

 Section 2 Design and Implementation

Figure 2=21 The Choice of a Segond-Best Path

raufing.algarithm.adamts.tp.the.abange.L.further...this...time

One can interpret the operation of the reachability elgorithm above as searching the network for augoessively longer paths ofter the best path fails, until it gives up at paths of length MAX. It is important to note that while conducting this search for a path of some length, the routing algorithm must continue to believe that the node is reachable, and thus must route traffic to it.

Now that we have discussed what makes the simple distributed reschability algorithm slow, and why these effects are common to a wide class of distributed algorithms, we present a partial ablution. It is disarly desirable to reduce the rise time of the routing elgorithm if that can be done without impairing its effectiveness. One way to think about this goal is to treat it as collapsing the greeby-one search for paths of increasing length into a single search for paths of any length, and choosing the best one. This highwievel idea still leaves open the problem of use as the best route while the seerch is in progress, while the ald best route is no longer good but no replacement. has One possibility, instead of oscillating between been found. verious routes as the search proceeds, is to continue to use the earlier best path until a new, better path is found. Now consider this technique for handling the transition period. Is it not lust the technique we need to speed the search (tself? Taking the value of the old best route as still the best has exactly the desired effect in the search. If there are any routes with bester values they will be found in the time it takes the network to perform a search, or what we termed above the lag time. That is, it takes only one step to find a better route if one exists. In the case that the old best route now is unusable that a length of MAXI, the node weits only the specified time period, the lag time, for a better route to be found. If a peth does not turn up in that time, it can be assumed that one does not exist. This is the partial solution we want, which we will deli hold downs

Ihairautiag.algarithmighmuld.agatimus.ta.uga.tha..bhat..raytg
to.a.givah.gastimatish..both.for.updatima.ang.faswarding..tas
anma.tima.ogrica..after.it.gath.worsa.

That is, the algorithm would report the surrent value of the previous best route to its neighbors and use it for routing packets, for a given time interval. This idea was first suggested by W. R. Crowther of BBN (internal memo, 1971).

Another way to look at this is to distinguish between changes in the network topology and traffic thet nesessitate changing the charge of the best route from these changes which merely affect the characteristics of the route, like hop count, delay, and throughput. The latter may represent a much larger class than the former, particularly in large networks, While the actual length of the path may change, the general direction in which packets

characteristics may.

should be routed is not likely to change, given that nodes are widely distributed and have only a few lines. This is what motivates believing the best line for a few tieks longer. identity of the best path probably does not change, though its

1. In the case when the identity remains the same, the mechanism of hold down provides a stepewise adeptation to the ghanges in the characteristics of the path, that is, a rise and fall time of sero, Certainly, this is optimal.

2. When the identity of the path must change, the rise and times are both equal to one las time. This is optimal for any algorithm within the prectical limits of propagation times discussed above.

enother interpretation of this technique is suggested by the network configurations of a chain of modes, such having 2 lines. In this network it is older that when the first mode, call it 0, open down an intermediate node, say node 2, thinks it has an alternate path, longer then its original path, via its other line, Of course this is not true, because the hop count at node 3 is besed on the hop count at node &, and hee no meaning if the best path at mode 2 is not longer usable. Clearly, this mituation holds in general because of the way that the distributed computation amouths the routing information over the network, Finds no nade holds a hop sount which differs by more than one from that held by any of the adjacent nodes, there will dertainly be appeared alternate paths out all its other lines. The technique of holding the original best path as still the best, with its new value, has the effect of <u>purging</u> all those sites depending on that path as a subspath for their routing. Within one led time, no nodes in the network will be using the old information. If the other lines at the node are reporting better paths than the original one, the node one than believe that they are new paths, independent of the original path. Thus a node enters a transitional state, which we call hold down, when it needs to adquire acquire information about a new best route.

There are two points which should be clarified about when hold down should be invoked.

The mode should ester held down when the value of the becomes line has the ester line becomes better the serious best line.

The reason for this distinction becomes diser when one considers the motivation for hold down laid out above. The purpose of the hold down mechanism is restricted to fording stepewise adaptation to changes for the worse in routing information. The algorithm should be allowed to adapt freely to imprevements in the routing information being received on any line, the previous best line or any other. A second clerification to the mechanism is as follows:

Clearly, there must be a procedure to ecquire a new route, and the whole dame for hold down is based on the abservation that the acquisition of the new reute cannot proceed reliably until all the (nformation has been purced. There are two practical considerations which should be mentioned. First, the hold down time should be run in parallel with the line-soing-down time rether then in serial, efter the line has some down. That is, if hald down rung for the lest a seconds of the lineacoing-down time. then a new route should be established (or the absence of any route established) by the time that the line is designed unusable, A second point to make is that this hold down meghanism can be invoked quite conveniently by introducing a dummy input routing message from the dying line, containing entries destinations set to MAX. This has the effect of dauging the algorithm to enter hold down for just those nodes for which the dving line was the best route.

Now we turn to the question of the duration of the hold down mechanism. For this purpose there must be a timer corresponding to the best path to each destination, slong with the identity of the best path, RH in the algorithms above, and the value of the variable being optimized, D above when we considered distance. The first point to make about this timer is:

If the best line date worse and them date better . so in . . the bald down times "nust be allowed to the to seemple tien.

A simple example suffices to illustrate that reacting the timer if the information from the best line improves may gause the original degradation in the line not to be communicated. Consider the following network, with daily of unity in both directions on both lines:

Figure 2-22 The Hold Down Timer Must Complete

Suppose the delay from node 2 to node 1 increases to 5, then decreases to 4, while other delays remain constant at 1. The delay to node 1 seen at nodes 2 and 3 is as follows: node 2 node 3

1 2

5 6 begin hold down

4 5

...

4 5 and hold down

Now suppose we slips the improvement in dalay from node 2 to node 1 to reset the hold down timer. It may happen that node 3 does not heer about the increase in delay to 5 until after it sends routing to node 2. This routing then may arrive at node 2 just after the improvement in delay from node 2 to node 1. When the

improvement is detected, hold down is terminated, and the information from node 3, out of date and incorrect, is believed. The delays can be shown as

node 2 node 3

- 1 2
- 5 2 bagin hold down at node 2
- 4 2 and hold down at node 2
- 3 2 incorrect information from node 3

edepted by node 2

This example suggests two more observations. The first concerns evolding the problem, described above, of introducing apprious routing information into the network.

Ite.held.daws.time.must.be.long.enade.for.the.couting information.te.organete.from the made.fin.cuestion.to..ell adlacent.bades.bad.back.agsin.

The key function of hold down is to purpe the old routing information from any neighboring nodes so that they do not egho it back and thus impede the rapid adaptation to new information. For this end it is not necessary to wait for the entire natwork to respond to the change in routing, but only for the immediately adjacent nodes. Actually, this statement may be too strong in some gases. For example, in the network in Figure 2-22 the best routes (besed on delay, say) to node 1 ere indicated by arrows, In the event that the line from node 3 to node 2 gets worse, node

3 will enter hold down. Then node 6 will soon initiate a hold down because its best line has become worse. Node 4, on the other hand, is not effected since its best route is not by way of node 3. It will not begin to hold down until nodes 3. 6. and 5 have all begun their hold down timers. Thus, node 3 must keep hold down on for long enough to allow the new information to propagate around the cycle. Otherwise, its hold down time will expire and node 4 will report old information that is no longer valid. The conclusion is

Ibe beid down time must be long enough for the routing information to brongers through several hops.

The exact number of hops depends on the likelihood of a path of several hops being superior to one of fewer hops. With minimum hop routing, hold down need last only long enough for routing information to propagate one hop. If the objective function is delay or throughput, it must in general be longer. The penalty for a hold down timer which is too short is not severe; routing will oscillate when the recently-acquired elternate route is found to be no better after all. No permanent incorrect information is introduced.

A second observation related to the natwork shown in Figure 2-23 is that if the line from node 2 to node 3 cerries routing updates less often then the other line, its hold down time must be longers

Figure 2=23

Temarinaria antinarianation and the belarantine and the the temperation and the temper

Ideally, the node has separate hold down timers for each line, running at independent retes, holding down the transmissions. In practice, it may be easier to adopt a set of hold down counters for each line speed present at that node.

The first point to make about reachability is one that we overlooked previously, but which essumes more importance now. Consider the problem of having the reachability algorithm adapt efter the network has been isolated into two or more components, as shown in Figure 2-24. Suppose the line between nodes 1 and 2, which is the only connection between Networks 1 and 2, has been down and now comes back up. Nodes 1 and 2 exchange routing information. Node 2 reports that all of Network 2 is up. Node 1 finds out about this, but it cannot initiate a message to a node in Network 2 because mone of those nodes have yet been informed that node 1 is up. Data packets can travel faster than routing information. If node I were to send a massage to one of the nodes in Network 2, the returning RFNM and any replies from the Host would be last because the destination node would think node 1 was still unreachable. Thus when node I sees all the nodes in Network 2 come up for the first time, it must mark them as in a special Chapter III Section 2 Design and Implementation

coming-up state. While a node has a given destination node marked as coming up, it can accept store-end-forward traffic for that node, and can propagate alive routing about that node, but it cannot initiate any Host traffic for it. The timer must be long enough for routing to propagate through the whole network (singe Network 1 may be a single node), though it can be signed as soon as a message is received from a node in the coming-up state.

While the need for a soming up delay is common to any reachability agheme, there is a similar delay which is necessary only because of the nature of hold down. A going down delay is elso reseasery, as can be demonstrated by the large ring shown in Figure 2=25. Suppose the shortest paths to node 1 from ell other nodes are indicated by the arrows and the line from node 1 to node 2 goes down. Then node 2 begins a hold down, reporting a value of MAX hops to node 1 in its routing message to node 3. Then node 3 in turn begins a hold down, and so on. Eventually, node X holds down but node Y keeps its old route. By this time, the hold down at nodes 2 and 3 may be over if the loop is large. They must now wait for the information about the other route to propagate from node Y to node X, and eventually to nodes 3 and 2. Thus the going down timer for a node must also be long enough for the network. to respond to a change in the routing. Periodically, in concert with the running of the shortest path algorithms the node makes the

Section 2 Design and Implementation

Figure 2=24 The Need for a Coming Up Delay per Node

Figure 2=25 The Need for a Going Down Delay per Node

determination as to whether all destinations are reachable on the basis of whether HOP(I) is less than MAX or not. And as we have seen, the transitions in state must be damped according to specified time constants. As shown in Figure 2m26, each destination node is in one of 4 states as seen by each of the other nodes in the network. The constants NCU, coming up time, and NGD, going down time, are topology-dependent. Here "reachable" refers to whether HCP(I)=MAX or not, and "get message" means that a message is received from the destination node.

The state diagram for the reschability of a given node resembles closely the atete diagram indicating whether a given line between nodes is useble or not. The logic shown in Figure 2-27 is from the line state program used in the ARPANET, but the principles are quite general. The IMPs decide whether a line is up or down on the basis of "Hello" and "I Heard You" (IHY) messages. There is a quality gentral placed on the line by the two constants LGD and LCU. If the IMP misses LGD I Heard You's in a row, it declares the line unuseble and weits for LHD (Line Held Dead) ticks so that the other IMP learns that the line has gone down. Then the IMP must redeive LCU IHY's before the line dan down. Then the IMP must redeive LCU IHY's before the line dan down. The line is reset if LR ticks go by without success.

Terminology. We begin by explaining the terminology used below to describe the ARPANET routing algorithm. The ectual IMP program is written in assembly language and does not use the terms

Section 2 Design and Implementation

Figure 2-26 Node State Diagram

Figure 2=27 Line State Diagram

employed here: we have reswritten it below in a PL/I=like The aldorithm Performs. two baste anguage. funations: determining whether each other IMP is reachable, and calculating the path of least delay to those IMPs which are reachable. The program maintains a number of tables indexed by IMP; these all have one entry for each of the NN IMPs in the network. All tables carry checksums which are a function of the NN entries, and which the program periodically verifies in order to minimize and localize the effects of any hardware or activare failures, the directory giving the output line corresponding to the shortest hop path to all IMPS; DROUTE is the analogous minimum delay directory. HOP and DEL are the tables containing the number of hops and estimated delay on these paths. The maximum values in these tables are denoted by HMAX and DMAX respectively. HOPIN and DELIN are the two parallel tables which together constitute the routing message received at an IMP, HOPOUT and DELQUT are the tables constructed by an IMP to form its next routing message to be transmitted, and to be used as the new HOP and DEL tables. elgorithm uses a technique seiled hold down to ensure that an IMP will not been its routing decisions on old information sent to lit its neighbors, The technique consists of using a timer, HOLDI for the best hop path and OHOLDT for the minimum delay path. wait a short interval before deciding to change to a new path. This interval is long enough for all the neighbors of the IMP

isern about the new information, so that the IMP sen make (ta decision addurately, rather than on the basis of "edhos" of its old information. This interval is set to HTMAX. In order to determine if hold down is negessary on the minimum delay path, the program keeps two dounters, DELOLD and DELTHS, which accumulate the indrease in dalay in the previous update period, and in the durrent update period respectively. Finally, the program inserts a transition delay so that all IMPs in the network learn about the transition before it begomes final, GODOWN is a table of timers which are set to GDMAX when an IMP begins to go down; COMEUP is a similar set of timers set to CUMAX when an IMP is first detected as reachable.

Routing Message Input Processing. When a routing update message is received, the input module first verifies the checksum on the message. If it is incorrect, the program reports an inter-IMP failure and ignores the message, Otherwise, it puts the message on a queue for a lower-priority routine. This routine then takes the following actions:

- 1. Take the routing message off the input quaue,
- 2. Verify the checksum on the following gode before executing it,
- 3. Abort if the checksum is incorrect.

- 4. Ignore the meseage if the input line is down, or if it is looped, or if the seriel number on the routing message is the same as that on the last message requived, or if the format compatibility number does not metal the number used by the raceiver.
- 5. Initialize the checksum on the output message to be generated, put in the sender's node number, compatibility number, and the ingremented serial number,
- 6. Perform the besid routing computation as follows:

```
For Idel to NN dos
  }f Im5gLP them HOPOUT(I),DELOUT(I),ROUTE(I),DROUTE(I)∢⇒0;
    elee dor
      IF ROUTE(1) WRMAX
        then /* I has been declared unreachable */
          1+ HOPIN(I)≥HMAX
            then /* I is still not reachable via J */
              do: DELOUT(J) <- DMAX; HOPOUT(I) <- HMAX;
              endi
            else /e a new path to I has been discovered a/
              do; COMEUP(I) <= CUMAX; /= Set Come=up time= =/
                  ROUTE(I) <= J+
                  HOPOUT(I) <= HOPIN(I)+1;
        else /# I has not been declared unreachable #/
          if ROUTE(I)aJ
            then /* Current minehop path is via J */
              dos (+ HOPIN(I)+1>HOP(I)
                     then /* Best path got worse, hold it down */
                       HDLDT(I) 46HTMAXs
                  IF HOPIN(I)+14HMAX
                     them /* Best path still exists */
                       do: GODOHN(I) 4+B;
                           /* Make sure Godown timer is off */
                           HOPOUT(I) 4 wHOPIN(I) + 1 y
                       ends
```

endi

message bufter.

Finalize the souting message sheeksum and put it in the

- 8. If the surrent routing message buffer is not being transmitted on any line, circulate the triple buffer system as follows:
 - 1. essign the new message to the output buffers
 - 2. Assign the output buffer as the free buffer areas
 - 3. essign the free buffer area as the buffer in which the next message will be built.

Routing Message Dutput Processing. When the progrem redelves the completion interrupt on one of the IMP's lines, it first checks to see if the senderouting flag is set. Routing transmissions take priority over all others. If the flag is set, the program takes the following steps:

- i. Verify the checksum on the following code before executing it.
- 2. Abort if the checksum is incorrect.
- 3. Locate the next routing message buffer from which to transmit, and verify that the message has a correct checksum;
- 4. Report an intrawIMP failure if the checksum is in error, and resynchronise the line protocol before continuing.
- 5. Otherwise, initiate the output and dismiss the interrupt,

Periodis Routing Processing. Every 648 militageonds the program takes the follow steps:

Section 2 Design and Implementation

- 1. Verify the checksum on the following code before executing it.
- 2. Abort if the checksum is incorrect.

```
For Idel to NN dos
  if HOLDT(1)>0
    then /4 degreement hold down timer for hope; 1,25=1,92 sees 4/
    HOLDT(I) 4=HOLDT(I)=1;
  if DHOLDY(I)>Ø
    then /* degreement hold down timer for delays 1,28=1,92 secs */
    DHOLDT(I) <- DHOLDT(I) - 1;
  DELOLD(I) a=DELTH8(I); DELTH8(I) 4=0;
  BELLINMODOS +1
    then /* decrement timer for IMPs going downs 7.7=11.5 secs */
    do: GODOWN(I) <=GODOWN(I)=1;
        if GODOWN(I) #Ø
           then /* IMP I just died */
             do: ROUTE(I) <- RMAX; DROUTE(I) <- RMAX;
            ends
    endi
  ( COMEUP(I) >B
    then /= degrement timer for IMPs coming up: 54.6-48.5 gads =/
    COMEUP(I) <- COMEUP(I)-1:
  /* update checksums on any tables with changed entries */
ends
```

3, Verify the checksums on the various tables: ROUTE, DROUTE, HOLDT, DHOLDT, GODOWN, COMEUP.

Every 25 milliseconds, the program takes the following steps:

is Set the senderouting fleg for the appropriate lines edgerding to their line speed and line utilization. The rule used is that the overhead due to routing messages should range from \$% of the line bendwidth on busy lines to 15% of the line bendwidth on files.

Section 2 Design and Implementation

2. Attempt to ofreujate the routing message buffers as in step 8 in the section above on Routing Message Input Processing.

2.4.3 Source IMP we Destination IMP

2,4,3,1 Basia Concepts

There has been considerable controversy during the last several years over whether of not a storemendatorward subnetwork of nodes should concern itself with endatored transmission procedures, Many workers feel that the subnetwork should be close to a pure pasket carrier with little concern for maintaining message order, for high levels of correct message delivery, for message buffering in the subnetwork, etc. Other workers, including the ARPANET designers have felt that the subnetwork should take responsibility for many of the endatorend message processing procedures. Of course, there are some workers who hold to positions in between. However, many design issues remain constant whether these functions are performed at host level or subnetwork level, and we discuss these constants in this section:

Buffering and Ripelising. As noted earlier, any practical network must allow multiple messages simultaneously in transit between the source and the destination, to achieve high throughput. If, for example, one message of 2000 bits is allowed to be outstanding between the source and destination at a time, and the normal network transit for the message including destination-to-source acknowledgment is 100 milliseconds, then the throughput rate that can be sustained in 20,000 bits per second.

Section 2 Design and Implementation

If slow lines, slow responsiveness of the destination host, great distance, etc., seuse the normal natwork transit time to be half a second, then the throughput rate is reduced to only 4,988 bits per second. Similarly, we think that pipelining is essential for most natworks to improve delay characteristics; data should travel in resconebly short peckets.

to summerize, low delay requirements drive packet size smaller, network and host lines faster, and network paths shorter (i.s., iewer nodestouned hope). High throughput requirements drive the number of packets in flight up, packet overheed down, and the number of elternative paths up.

<u>tror. Camtral.</u> We consider source-to-destination error control to comprise three tasks; detecting bit errors in the delivered messages, detecting missing messages or pieces of messages, and detecting duplicate messages or pieces of messages.

The former task is done in a straightforward manner through the use of checksums. A checksum is appended to the message at the source and the checksum is checked at the destination; when the checksum does not check at the destination, the incorrect massage is discarded, requiring it to be retrememitted from the source. Several points about the manner in which checksumming should be done are worthy of note. (a) If possible, the checksum should sheek the correctness of the resequencing of the messages

which possibly got out of order in their traversal of the network, (b) A powerful checksum is more efficient than alternative methods such as replication of a critical control field, it is better to extend the checksum by the number of bits that would have been used in the redundant field, (g) Unless encryption is desirable for some other reason it is simpler (and just as safe) to prevent delivery of a message to an incorrect host through the use of a powerful checksum than it is to use an encryption mechanism; (d) Node-to-node shecksums do not fulfill the same function as end-to-end checksums because they sheck only the lines, not the nodes,

An inherent dheresteristic of packetsswitching networks is that some messages or portions of messages (i.e., packets) will fail to be delivered, and there will be some duplicate delivery of messages or portions of messages, as described in the section on network properties. (Please note that throughout the remainder of this subsection we use the word "message" to mean either messages or portions of messages (i.e., packets).)

Missing messages can be detected at the deatination through the use of one state bit for each unit of information which can be simultaneously traversing the network. An interesting detail is that for the purposes of missing message detection, the state bits used must precisely cycle through all possible states. For

example, stamping messages with a time stamp does nothing for the process of missing message detection because, unless a message is sent for every "tick" of the time stemp, there is no way distinguish the case of a missing message from the case where no messages were sent for a time.

detected with an identifying Duplicate messages can be sequence number such that messages which afrive from a prior point in the sequence are recognized as duplicates. What should be noted carefully here is that dupligate massages can arrive at the destination up to some time, possibly quite long, after the original copy, and the sequence number must not complete a full eyele during this period. For example, if a network goal is to be able to transmit 200 minimum length messages per second from the source to the destination and seah needs a unique sequence number, and if it is possible for messages to arrive at the destination up to 15 seconds after initial transmission from the source, then the sequence number must be able to uniquely identify at least peckets. It is usually no trouble to deloulate the maximum number messages that can be sent during some time interval. What is more difficult is to limit the maximum time after which duplicate messages will no longer arrive at the destination. One method is to put a timer in each message which is counted down as meseage traverses the network; if the timer ever counts out, the message is discorded as too old, thus quaranteeing that no messages older than the initial setting of the timer will be delivered to the destination. Alternatively, in practice one can make a resecrably good approximate deliquiation of the maximum arrival time through study of all the worst case paths through the network and all the worst case combinations of events which might cause messages to loop around in the network.

wither case, there certainly must be mechanisms to resynchronize the sequence numbers between the source and the destination at node start-up time, to recover from a node failure. eta. A good practice is to resynchronize the sequence numbers occasionally even though they are not known to be out of step. A good frequency with which to do redundant resynchronization would be every time a message has not been sent for longer than the maximum delivery time. In fact, this is the maximum frequency with which the resynchronization can be done (without additional meghan(sms); if duplicates are to be detected reliably, the sequence number at the destination MUSE function without disruption for the maximum delivery time efter the "lest messade" has been sent. If it is desirable or necessary to resynchronise the sequence numbers more often than the maximum time, an additional Husek number must be attached to the sequence number to uniquely identify which Pinstance" of this set of sequence numbers is in affect; and, of course, the packets must also carry the use number.

The next point to make about endetowend error dontrol is that any message going from source to destination can potentially be missing or duplicated; i.e., not only data messages but control messages. In fact, the very messages used in error control (e.g., sequence number resynchronization messages) can themselves be missing or duplicated, and a proper endetowend protocol must handle these cases.

Finally, there must be some inquiry-response system from the source to the destination to complete the process of detecting lost messages. When the proper reply or acknowledgment has not been required for too long, the source may inquire whether the destination has received the message in question. Alternatively, the source may simply retransmit the message in question. In any dage; this source inquiry and retrensmission system must also function in the face of duplicated or lost inquiries and inquiry response control messages. As with the interanode acknowledgment and retrememission system, the end-to-end acknowledgment and retransmission system must depend on positive acknowledgments from the destination to the source and on explicit inquiries or retransmissions from the source. Negative acknowledgments from destination to the source are never sufficient (because they might get lost) and are only useful (albeit sometimes very useful) for increased efficiency.

Storage Alicestics and Flow Control: One of the fundamental rules of sommunications systems is that the source cannot simply send data to the destination without some mechanism for quaranteeing atorage for that data. In very primitive systems one querentee a rate of disposal of data, as to a line printer, and not exceed that rate at the data source, In sophisticated systems there seem to be only two alternatives, Either one can expidditly reserve space at the destination for a known amount of data in advance of its transmission, or one can desiare the transmitted copy of the data expendable, additional copies from the source until there is an acknowledgment from the destinction. The first elternative is the high bandwidth solution: When there is no space, only tiny messages travel back and forth between the source and destination for the purpose of reserving destination storage. The second elternative is the low delay colutions the text of the message propagates as fast possible.

In either case storage is tied up for an amount of time equal to at least the round trip time. This is a fundamental result we the minimum amount of buffering required by a communications system, either at the source or at the destination, equals the product of round trip time and the channel bendwidth. The only way to circumvent this result is to count on the destination behaving in some predictable fashion (an unrealistic assumption in the general case of autonomous communicating entities).

As we stated earlier, our experience and enalysis convinces that if both low delay and high throughput are desired, then there must be mechanisms to handle each, since high throughput and low delay are conflicting goals. This is true, in particular, for the storage allogation mechanism. It has occasionally been suggested, mainly for the make of simplicity, that only the low delay solution be used; that is, messages are transmitted from the source without reservation of space at the destination. Those people making the choice never to reserve apage at the destination frequently assert that high bandwidth will still be possible through use of a mechanism whereby the source sends messages toward the destination, notes the arrival of acknowledgments from the destination, uses these acknowledgments to estimate the destinction reception rate, and adjusts its transmissions to match that rate. We feel that such schemes may be quite difficult to parameterize for efficient control and therefore may result in reduced effective bandwidth and increased effective delay. If the source never sends to the destination so fast that the destination must discard enything, then the delay is very low, but the throughput is not as high as it might be. Further, unless the source pushes now and then, it will never discover that the destinction is able to increase its throughputs. On the other hand, when the source is pushing hard snough, the destination may suddenly cut back on its throughput, sausing all the messages

which will be discarded at the destination due to the sudden dutback to have to be retransmitted, increasing effective delay. the destination could be predicted to accept traffic at a steady rate and wary this rate only very mlowiv. the type of feedback avstem aited above might work. In this case, unacknowledged messages should be retrensmitted from the source to destination shortly after the expected time for the acknowledgment to return has elapsed, if minimum delay and maximum throughput are to be obtained (this is in contrast to the often suggested practice of keying retransmissions to the discard rate), However, in predtice, the time for the acknowledgment to return is likely to be very difficult to predict due to veristions (possibly rapid) in the transit time of the communications channel and particulariy in the response time of the destination. Furthermore, the greater the sum of transit time and response time, the looser and less efficient the feedback loop will be. In there appear to be oscillatory conditions which can occur where performance degrades completely. (Note that if there fo much possibility of message loss, then the acknowledgment and retrensmission system should allow quite selective retransmission of measages rather than, for (netange, requiring a complete window ΔÌ messages to be retransmitted to affect retransmission of the specific messages requiring its otherwise, message retransmission will use expessive bandwidth.)

abaye . discussion assumes that all mechanisms are attempting to minimize the probability of message discard. If, in possible discards addition at the destination the to communications chennel solves its internal problems . potential deadlocks) with cavaller discarding of messages, or if destination actives its internal problems with daysliar discarding of messages, the detrimental effects of discarding (reduced effective bandwidth and ingressed effective delay) are probably drastically increased. Further, the above discussion assumed the destination was able to minimize the probability of discard. While this may be possible for a single source, we think it is unlikely that the destination will be able to resolve, in a way that does not entail excessive dispards, the contention for destination storage Pom multiple ungoord nated sauroes. Detrimental contention for destination storage, in the absence of a storage reservation mechanism, happens practically continuously under even modest treffic loads, and in a way uncoordinated with the rates and strategies of the various sources. As a result, well-behaved hosts may unavoidably be penalized for the actions of poorly=behaved hosts.

In addition to space to hold all data, there must also be space to record what needs to be sent and what has been sent. If a message will result in a response, there must be space to hold the response; and once a response has been sent, the information

about what kind of answer was sent must be kept for as long as retransmission of that response may be necessary.

Presedence and preemption. The first point to note about precedence and preemption is that the total transit time being apecified for most packet-switching networks of which we are aware is on the order of less than a few seconds (often only a fraction of a second). Thus, the traditional specifications (for example, low priority treffic must be able to preempt all other traffic so that it can traverse the network in under two minutes) no longer make much sense. When all messages traverse the network in less than a few seconds, there is generally no need to specify that top priority traffic must preempt other traffic, nor to specify the relative precedences between the other types of traffic. Priority can be used, however, to admit traffic into the network selectively.

Though priority is not strictly necessary for speed, it may be useful for contention resolution. It appears to us that there are three precedence and preemption strategies that are reasonable to consider for a pasket=switching network, Strategy I is to permanently assign the resources necessary to handle high priority traffic; this guarantees the delivery time for the high priority traffic but is expensive and should only be done for limited high priority traffic. Strategy 2 is to preempt resources as necessary

DRAFT

for high priority traffic. This can have two effects. Preempting packet buffers results in data loss: Preempting internal node tables (e.g., the tables associated with packet sequence numbering) results in state information loss. State information loss meens that data errors are possible which may go unreported, Strategy 3 is not to preempt resources, and to rely on the standard machanisms with a priority ordering. This is simple for the nodes, but it does not of itself guarantee delivery within a certain time.

We think the correct strategy is probably a mixture of the strategies above, Possibly some resources, on a very limited basis. should be reserved for the tiny amount of flash traffic. This guarantees minimum delay without any quausing latenty. For the rest of the treffic, the normal delivery times are probably acceptable. The presence of higher priority traffic can cause gradual throttling of lower priority traffic, without loss of state information. As the time to do this graceful throttling is normally only a frection of a mesond, the higher priority traffic has no real reason to demand (natantaneous, information-losing preemption of the lower priority traffic.

Manage 51re. The question is often asked, "If one increases packet size and decreases message size until the two become the same, will not the difficult message reassembly problem be

removed?" The enswer is that, perhaps unfortunately, message size and packet size are almost unrelated to reassembly.

We have already noted the relationship between delay and packet size. Delay for a smell priority message is, to first order, propertional to the packet size of the other traffic in the network. Thus, smell packets are desirable, Larger packets become desirable only when lines become so long or fast that propagation delay is larger than transmission time.

Message size needs to be large because the overhead on messages is significant. It is inefficient for the nodes to have to address too many messages and it may be inefficient for hosts to have too many message interrupts. The upper limit on message size is what can conveniently be reessambled, given node storage and network delays.

When a ghannel has an appreciable delay, it is necessary to buffer several pieces of data in the channel at one time in order to obtain full utilization of the channel. It makes little difference whether these pieces are called packets which must be reassembled or messages which must be delivered in order.

We do not feel that the choice between aingles and multispecket messages is as important as all the controversy on the subject would lead one to believe. There is agreement that

buffering many data units in transit through the natwork simultaneously is a necessity. Having multispacket messages is probably more efficient (as the extra level of hierarchy allows overhead functions to be applied at the correct, i.e., most efficient, level); having single-packet messages probably offers the opportunity for finer grained storage allocation and flow control mechanisms.

Multiplexing and Addressing. Up to this point in our paper, have not been very specific about whether the above-mentioned flow control, sequencing, error control, etc. mechanisms performed for each pair of communicating processes, or whether several processes communicating between a given pair of source and destination nodes share a set of these control mechanisms. tradeoff is between overhead and presision of control. If many conversations. 870 multiplexed on eadh instance 01 source-to-destination control mechanism, the control overhead is lower than if each conversation has its own control mechanism. 0 n the other hand, if several conversations are multiplexed on the semm control mechanism, all the conversations tend to have to be treated equally (e.g., if one is stopped, all are stopped); while 10 each conversation has its own control mechanism, decisions about the allogation of various resources to the various conversations can be made. To give some examples of the latter, conversetions over separate control mechanisms can be given

differing allocations, priorities, treatments of error conditions, etc.

Another issue is the management of the space available for control mechanisms when it is insufficient to handle the number of conversations competing for the communications channel. late-comera be left out until resources are evallable, or should some way be found to multiplex the available gontrol mechanisms in time among the demanding conversations? We believe the latter should be done. The key here is not to allow users to explicitily sequire and hold resources (e.g., control mechanism space) for interprocess communication. Instead, the system should notice which users are actively communicating and dynamically gather the needed resources by garbage=collecting the resources previously being by users which appear inestive. This dynamic ueed essignment of resources is obviously not fundamentally different from the acheduling of eny limited resource in an operating system (e.d., memory, CPU cycles, the I/O channel to the disk; and therefore has all the mormal possibilities thrathing, forunfalaness, and so on, if care is not taken.

Once decisions in all of the above areas of multiplexing are made, one must choose the addressing machanism and formats to be used. This is usually quite straightforward. The main point here is that addressing comes last; but very often we see designs

begun by choosing the addressing system and format: A similar statement can be made about the choice of all other message formats:

2,4,3,2 Original Implementation

In this section we describe the mechanisms by which the IMPs in the ARPANET manage the flow of data from a source host to a destination host.

2.4.3.2.1 Message Format

In the ARPANET, a host presents messages to the IMP to which the host is directly connected. These messages must be less than about 5190 bits long, and the messages are transmitted between the host and IMP over the host/IMP interface, This interface has two parts, the hardwere part and the software part.

The hardware part of the host/IMP interface (tself has two parts, a standard portion supplied with the IMP which is (almost) identical for all hosts and a special portion supplied by the host. For simplicity and power, the standard interface has been defined to be full duplex, Electrically, the standard interface follows a bitebyebit handshaking procedure. The procedure is something like "I'm ready to send a bit," "I'm ready to receive a bit, "" "Here's the bit, "" "Good; "" The procedure also provides for saying "That's the last bit in the message, "" In

our opinion this procedure is important. The ARPANET has to connect together all kinds of different computers with different word langths, different speeds, different loading, and so forth; and it is desirable to place only minimal constraints on the hosts behavior. The above procedure permits this, The asynchronous, bit=by=bit serial interface permits both the host and the IMP to be able to start and stop transmission whenever necessary. Incidentally, the IMP's side of the above procedure is implemented entirely with hardware.

An optional synchronous host/IMP interface is also available. It is considerably more complicated, less flexible, and slower than the normal asynchronous, serial interfaces and its use is recommended only when an communications interface suitable for operation over long distances is required. HDLC would have been a better choice had it been sveilable at the time.

The software interface between an IMP and a host is also simple, using a minimum number of control messages. The host specifies to its IMP the destination of a message and a few other things in the first 32 bits of the message, salled the leader. Hessages arriving at the hosts have the same information in the first 32 bits of the messages, except that the destination is replaced by the source.

DRAFT

Neither the hardware nor the zoftware interface between the IMP and the host places any constraint on the content of messages other than that they must have legal leaders and must have less than the maximum length. In other words, mestages may be sent through the network containing arbitrary sequences of bits.

An IMP breaks messeges arriving from its hosts into pagkets 1000 or fewer bits long. As the IMP segments a message into packets, it appends to the front of each pecket some control information salled the header, The header contains destination, the packet number within the message, a message sequence number which is used for reconstructing the message stream as the messages arrive at the destination, and other control information (e.g., priority information) copied from the message into each constituent packet. This message segmentation and message reconstruction is completely invisible to the heats in the ARPANET implementation of messages and packets.

Sourcesto=Destination Transmission 2.4.3.2.2 The Ortainel Procedure

Having given this generalized introduction to the subject of measage processing, we turn to an examination of the original ARPA network implementation, in effect from 1975 to 1972. This system based on the concept of a link, a logical connection between hosts, on which only one message at a time could be two

transmitted. The source and destination IMPs kept tables of active links over which messages had recently been transmitted. This technique performed several of the functions described above, though some in a rather rudimentary fashion. As we shall see, some of the functions noted above for the facureer and idestination were performed in the source IMP and destination IMP, while others were left pertly in the domain of the source host and destination host.

First, the buffering in the network was up to the host program; it dould send only one message per link but the IMP placed no limits (in prestice, a very high limit) on the number of links that could be used. Pipelining was provided for, and still is, by breaking up messages (up to 5000 bits0 into packets (up to 1888 bits) to decrease delays. A reasonbly process at the destination is necessary to perform ordering, and is also implied by pipelining. Error control was performed by both the source and destination IMP by maintaining a sequence number to the messages on each active link. Each message and RFNM (Ready for Next Message, the delivery confirmation) darried this (dentifying number, and thus missing and duplicated transmissions gould easily be detected. Message sequenting was no problem with the one-message-on-a-link rule, since the destination host was then quaranteed to receive an ordered data stream on each link being used by the source host. Storage allocation was quite primitive;

the source IMP did not hold a copy of any message, nor was storage at the destination IMP reserved. Instead, the packets of the message were simply transmitted to the destination, where they were accepted if there was sufficient storage, or otherwise rejected; that is, the destination IMP would not send an ack back to the previous IMP, which would subsequently time retrensmit the packet. This approach therefore used the of the storemend-forward subnetwork when the destination IMP was full. Finally, flow control was likewise sudimentary and also somewhat in the province of the hosts, since it rested primarily on the technique of parmitting only one message per link, concept provided flow control to the extent that each link was well-controlled; the rate at which the source host could send was tied to the rate of the destination host, on a given link. But if too many links were used, no adequate flow control measures were present.

A brief critique of this first approach can be offered: major advantage, to our perspective new, was its eimplicity and therefore its basis reliability. The link table in the source and destination IMPs was dynamic, addressed by a heat key, and garbage-sollected on a time-out. The format of the table is illustrated in Figure 2=26. The pre-message rule meant that the link tables were easy to deal with and were solfesynchronizing, It was fair since all hosts had equal access to the common line

Figure 2028 Link Table Format

tales and storage pool in the IMPp. This technique was also efficient in terms of ellowing hosts great freedom to echieve minimum delays and meximum throughput levels, since it did not impose any overhead transactions. It was not efficient in IMP storage, since a fairly large has table entry was needed to despribe a single message, or in IMP bandwidth, since the hashed access was time-consuming.

The major oriticism of this askame (a s(mp)e) it can lead to a storedeebased deadlock that we sail reassambly lock-up, as illustrated in Figure 2-29. This deadlock occurs under high load, which first gauses reassambly congestion, in which efficiency is reduced due to retransmissions.

In Figure 2029, IMP 1 is sending multipacket messages to IMP 3; a lockwap can again when all the ressemble buffers in IMP 3 are devoted to pertially ressembled messages A and 8; Since IMP 3 has reserved all its remaining space for awaited packets of these partially ressembled messages, it can only take in those particular packets from IMP 2; Those outstanding packets; however, are two hops away in IMP 1; They cannot get through because IMP 2 is filled with storesendsforward packets of messages C, D, and E (destined for IMP 3) which IMP 3 cannot yet accept; Thus, IMP 3 will never be able to complete the resembly of messages A and B.

Figure 2-29 Reassembly lock-up

We were always aware that hosts could defeat this flow dentral mechanism by *apraying* messages over an inordinately large number of jinks, but we counted on the nonemaligious behavior of the hosts to keep the number of links in use below the However, simulations and which problems occur. level at experiments artificially loading the network demonstrated that communication between a pair of hosts on even a modest number of links could defeat our flow control mechanisms further, it sould be defeated by a number of hosts communicating with a common site even though each host used only one link, Simulations showed that reassembly lock-up may eventually oddur when over five links to a perticular host are simultaneously in use. We should stress that ordering is the besis cause of the problem. not reassembly of multi-packet messages.

Note that much of the trouble with this approach to message processing is that the responsibility for tasks such as buffering, sequencing, and flow control was shared between the IMPs and the hosts. That is, the IMP was responsible for sequencing, flow control, etc., with regard to a single link, but the host was responsible for the number of links that were active at once. As we shall see, we changed the ARPA network so that the IMPs had complete and effective sontrol over all the message processing functions in the communications subnetwork, while still leaving the hosts freedom to exercise their own controls as well,

2.4.3.2.3 The Next Implementation

In midel 972, the message programing algorithms just described were changed in several substantial ways. First, reassambly lock-ups were prevented by providing explicit storage allocation. Secondly, the sequencing and error control schemes were modified to be more general and more efficient. These techniques have remained in the ARPA network until late 1974. Further changes are described below in Section 2,4,3,3.

As we noted before, if storage is allocated one can optimize for low-delay transmissions. If the storage is allocated at the destination, one can optimize for high-bandwidth transmissions. To be densistent with our view of a balanced communications system, we have developed an approach to reassembly concention and lock-up that utilizes some buffer storage at both the source and destination; our solution also utilizes a request mechanism from source IMP to destination IMP.

Specifically, no multi-packet message is allowed to enter the network until storage for the message has been allocated at the destination IMP. As soon as the source IMP takes in the first packet of a multi-packet message, it sends a small control message to the destination IMP requesting that reassembly storage be reserved at the destination for this message. It does not take in further packets from the host until it receives an allocation

message in raply: The destination IMP quayes the request and sends the allocation message to the source IMP when enough reassembly storage is free; at this point, the source IMP sends the message to the destination.

Effective bendwidth is maximized for sequences of long massages by permitting all but the first message to bysoss the request mechanism. When the message itself arrives at the destination IMP is about to return the ReadysForeNext-Message (RFNM), and end-to-end delivery confirmation message, the destination IMP waits until it has room for an additional multi-packet message. It then piggy-backs a storage silocation on the RFNM. It the source host is prompt in enswering the RFNM with its next message, al eliocation is ready and the message can be transmitted at once. If the source host delays too long, or if the data transfer is complete, the source IMP returns the unused ellocation to the destination. With this mechanism, we have minimized the interemessage delay and the hosts can obtain the full bendwidth of the network.

The delay for a short message is minimized by transmitting it to the destination immediately while keeping a copy in the source IMP. If there is space at the destination, it is accepted and passed on to a host and an RFNM is returned; the source IMP discards the message when it regaines the RFNM. If not, the

message is disperded, a request for allocation is queued and, when space becomes available, the source IMP is notified that the message may now be retransmitted. Thus, no sateup delay is indurred when storage is available at the destination.

The above mechanisms make the IMP network much less sensitive to uneresponsive hosts, since the source host is effectively held to a transmission rated equal to the reception rate of the destination host. Further, reassembly lockeup is prevented because the destination IMP will never have to turn away a multi-packet message destined for one of its hosts, since reassembly storage has been allocated for each such message in the network. The east of this system is small; the one-time cost of implementing the more complex algorithm, plus a negligible overhead in line bendwidth.

Since the link mechanism was no longer needed for flow control (and had not been successful anyway), we felt that a less coatly mechanism should be employed for sequence control. The link mechanism was therefore aliminated from the IMP subnetwork. RFNMs are still returned to the source hosts on a message basis, but they are used only as asknowledgments to messages. To replace the permitted sequence control mechanism, we decided upon a sequence control mechanism based on a single logical pripate between each source and destination IMP. Each IMP maintains and

Section 2 Design and Implementation

independent message number sequence for each pipe. A sequence number is assigned to each message at the source IMP and this message number is cheaked at the destination IMP. All hosts at the source and destination IMPs share this message space. Out of an 8-bit message number space (large enough to accommodate the mettling time of the network), both the source and destination keep a amail window of gurrently valid message numbers, which ellows several messages to be in the pipe simultaneously, permitting high throughput rates. Messages arriving at a destination IMP with outsoference message numbers are dupliquies to be discarded. The window is presently four numbers wide: this perameter must be chosen in relation to the delay and bandwidth of the average and-te-end network path. The message number the two purposes of sequencing and error control.

Beduende dentrol svetom based On. single source/destination pipe, however, does not permit priority traffic to go aheed of other traffic. We solved this problem by permitting two pipes between each source and destination; a priority (or low-delay) and a non-priority (high-bandwidth) pipe, with duplicate detection performed in common.

hosts may, if they choosed, have several messages outstanding simultaneously to a given destination but, since priority messages fleepings shead, and the last message in a sequence of long

messages may be short, priority can no longer be assigned strictly on the basis of message length. Therefore, hosts must explicitly indicate whether a message has priority or not.

Hith message numbers and reserved storage to be accurately accounted for, cleaning up in the event of a lost message must be done derefully. The source IMP keeps track of all messages for which a RFNM has not yet been received. If the RFNM is not redelived within a certain time (presently about 30 seconds), the source IMP sends a control message to the destination enquiring about the possibility of the message been lost. The destination responds to this message by indicating whether the message in question was previously received o not. The source IMP continues enquiring until it receives a response. This technique guarantees that the source and destination IMPs keep their message number sequences synchronized and that any allocated space will be released in the rere dase that a message is lost in the subnetwork because of a machine fellure.

The format of the durrent message tables is shown in Figure 2=30. TMESS is the next message number to send; there are four bits to indicate whether the messages to whigh they refer (TMESS, TMESS=1, TMESS=2, TMESS=3) are still outstanding (unanswered). The ORDer field gives the next 2-bit ordering number to use on priority messages. RMESS and the essociated fields work in an

Figure 2-30 Message Table Format

analogous fashion for received messages, which are marked complete when they are given to the IMP=to=host routine, AMESS and RALLY are used to send back replies to messages in the window; AMESS is the next message to reply to; RALLY holds what kind of a reply to send. INC times out overdue transmissions; RST=T and RST=R time out efter a message number reset (as explained later) is performed on the transmit or receive side, to prevent more than one reset per 30 seconds, and thus detect duplicates; and RST indicates that a transmit reset is still in progress.

Several expects of this data structure are worth noting, first, the tables are quite compact compared to the sartier link tables, since a 4-word entry describes the complete state of four messages in flight between two IMPs, as opposed to three words in each of the source and destination IMPs to describe one message, Also, it is easy to access since it is indexed by IMP. Secondly, the structure of the RALLY table is important, both from the point of view of deedlock prevention and as a compact representation.

The destination IMP quoues up replies (such as RFNMs, ALLOCATES, etc) before they can be sent, An entry in this table is either blank or contains an RFNM, an ALLOCATE, etc. The IMP number to send it to and the message number it refers to are implicit in the position of the entry on the table. That is, the IMP has table storage to save all four possible replies that could

be pending, given the window size of 4, to all possible source IMPs (currently 64). This is important; the destination IMP siways has a piece to put a reply. When the reply is sent, the entry is marked as blank, and the number indicating which message number to enswer next, AMESS, is incremented.

Without table storage for each possible reply, in highly bit-goded form, the IMP would have to deal with queues of pending replies that dould grown very long. Of dourse, no limit can be plead on such a queue without running the risk of dauging resessably lock-up. All packets must be addented by the destination IMP from the storemend-forward part of the network to svoid such deadlock states.

In fact, the RALLY table is even more compact than it appears, since it can also keep a record of what kind of reply was sent back for each message number. That is, an entry for message N in the RALLY table is either a pending reply, or, if it is marked blank, a record of the reply sent for message N=4. Then, should the source IMP fail to receive the enswer and send an incomplete query 36 seconds later, the destination IMP can send back the correct original enswer. Note that no more entries need to be kept to do this. If RALLY holds enswers for messages 7, 5, 9, and 10, and then 7 is sent off, the destination IMP needs to remember the enswer is sent to 7 only until it has an enswer

queued up for message it. When it overwrites the answer for message 7, and the entry is marked full again, the source IMP will not now enquire about message 7; (to message number window has advanged pest that point.

A final appear of this structure worth noting is the resevonthronization machenism we have installed (we found that, in practice, all subnetwork systems need to be checked for consistency and fixed if wrong). It is suggested by the idea just noted the source IMP will legitimetely enquire about messages only up to one window's worth behind the gurrent position. This is best shown by explaining the actions the IMP takes under various direumstances, as illustrated in Figure 2011.

Case 4 is the one just discussed and case 5 triogers message number resynchronization. It says that if an IMP reserves an incomplete inquiry for a message outside of the range that the destination expects the source IMP to be using, it sends back a special outseferance message. Note that the offending incomplete may simply be an old duplicate. Therefore, the source IMP, when it receives the outseferance message, will sheek its message number to see if it if still weiting for the enawer to that inquiry. If it is, a fmessage number look-upf exists (presumably, this results from some software or hardware failure affecting the message tables). If not, an old duplicate incomplete has landed

Cese	Action	Message Type	Position Relative to destination window
1	Accept	Not incomplete	In window
2	Discard as duplicate	Not incomplete	Out of window
3	Discard any packets of messa age in progress	incomplete	In window
4	Send beck duplie cate correct reply	Incomplete	In 4 numbers just prevetous to duprent window
5	Send bask en out-of-range message	Incomplete	Outside of window of He and also window of H previous

Figure 2=31 Astion by Destination IMP on Various Messages

at the destination, and the out-oference message can be ignored,

Now, in the case that a look-up has been detected, the IMPs take the steps shown in Figure 2-32 to free it.

This set of mechanisms allows the IMPs to:

keep more addurate measage addounting;

detect lock-ups with a single incomplete transmission time-out delays

report the frequency of certain events to the Network Control Center, includings

incomplete queries sent
out of range incompletes received
receiver resets

sender resets

(by compering the last two ARPANET operations personnel can try to understand the gauses for any message number look-ups that do occur, and perhaps fix the failing hardware or software.

In short, the system is event-driven, giving a maximum amount of diagnostic information while still performing the necessary function at relatively low cost in delay, complexity, or bandwidth.

ARPANET Completion Report Chapter III

DRAFT
Section 2 Design and Implementation

Source IMP

Destination IMP

Mark frametwineprograms bit, prohibiting traffic to this IMP.

Send a reset message.

get reset message.
Reset all requive tables.

Send back *reset reply* message

Receive reset reply message.

If freset in progress oft
is on, reset all transmit
tables, return an Incomplete
Transmission for all messages in progress to that

IMP, Clear the freset in progress bit.

Resume normal operations.

Do not believe any other reset message foo another 38 seconds, to ignore duplicate is more than 30 seconds oid).

Figure 2-32 Action by Source and Destination IMPs on Detection of Lock-up

2,4,3.3 Subsequent Modifications

The final topic to consider under the general heading of problems with the next ARPA network approach is that of adapting the various techniques described for use in large networks with hundreds of IMPs and thousands of hosts. The key point is that the cost of keeping linear tables indexed by IMF, host, or anything else, becomes prohibitive at some metwork size. The simplicity and reliability of such tables, and their resulting freedom from deedlock, are naturally desirable but sennot be meintained beyond a certain level because of the high cost. 1972-1974 tables were structured for a network of up to 64 IMPS (each IMP has 16K of memory), your hosts per IMP (plus 4 fake hosts internal to each IMP), with line bendwidth of 58Kbs typically, and round-trip delays of roughly 1/2-second. With the growth of the ARPA network, and the introduction of new technology in IMPs and dirouits, including satellites, all the parameters above must be reeskamined.

2.4.3.3.1 Dynamic Slocks

The approach we decided to take is to make all the data structures in the IMP esseciated with message progessing take the form of dynamic blocks, each containing a few words of storage. One example is the transaction block, which the source IMP creates when a message is initiated. This block keeps track of whether

the message has a copy kept at the source, whether it needs an allocate, and so on. A key function of the block is to hold a copy of the message header to identify the message. When the RENM returns from the destination, a RENM for the source host can be formatted in place in the transaction block, solving a difficult storage allocation problem. Thus transaction blocks can be on a host queue (RENM ready to be sent). In addition, they can hold the information about whether an ALLOCATE was sent back with the RENM.

Another example of the dynamic block is the reassembly block, which is kept at the destination IMP for the purpose of ordering the packets of a multi-packet message. It also contains the message header, and a list of which packets have arrived. If none has arrived yet, it constitutes a record of an allocation that has been sent to a particular host and IMP. It can also be on several queuess a free list, an active list, or a host queue (this last is an implementation option we did not choose).

In keeping with this general philosophy of table structure, we redesigned the current message tables to use massage blocks instead. These blocks are made dynamic by keeping them only as long as a conversation between two hosts is sative. This concept allows the set of hosts at an IMP to send messages to an arbitrary number of other hosts over a wide range of addresses, with a

limited number of message blocks. The message number tables on each block can work exactly as they do now, and the other message processing functions such as flow gontrol can also proceed accordingly.

In order to successfully and efficiently handle the large number of conversations with various hosts that can be simultaneously in progress, all of the date structures in the IMP associated with message processing take the form of blocks dynamically gethered from a pool of blocks, each containing a few words of storage. The alternative of keeping linear tables indexed by IMP, host, or anything else is prohibitively expensive as the network becomes large.

One exemple of such a dynamic data structure is the trensection block which the source IMP creates when a message is initiated. This block keeps track of whether the message has a copy kept at the source, whether it needs an allocate, and so on. A key function of the transaction block is to hold a sopy of the message header to identify the message. When a Ready For Next Message returns from the destination, a Ready For Next Message for the source host can be formatted in place in the transaction block, solving a difficult storage allocation problem. The transaction block, solving a difficult storage allocation problem. The transaction blocks can be on a free list, on an active list (i.e., message outstanding), or on a host queue (i.e., Ready for Next

Message to be zent). In addition, they can hold the information as to whather an "allocate" was sent back with the Ready For Next Message.

Another example of a dynamic block is the reassembly block which is kept at the destination IMP for the purpose of ordering the packets of a multi-peaket message. It also contains the message header and a list of which packets have arrived. If none have arrived yet, it constitutes a record of an allocation that has been sent to a particular host or IMP. It can also be on several queues; a free list, an active list, or (theoretically) a host queue (this last is an implementation option which was not chosen in the case of the IMP system). Notice that in each of these cases, that of the transaction block and the reassembly block, the date structure represents an important resource. There can only be a finite number of blocks; their use must be allocated among several competing source and destination hosts; critical questions of efficiency and fairness arise in the process of block allocation.

In keeping with this general philosophy of table structure, message blocks are used to keep track of the host/host message numbers. This concept allows the set of hosts on an IMP to send messages to an arbitrary number of other hosts over a wide range of addresses, with a limited number of message blocks. Several issues trian with consideration to dynamic message blocks:

There must be control messages between the source and destination IMPs of the form "get a block" and "got a block"; in order to establish a "conversation" between a given pair of hosts on the two IMPs.

There must be an error control mechanism to detect duplicate or missing aget a block and aget a block messages.

Once a conversation is established, massages can flow. Then there must be a technique to distinguish messages in this conversation from old duplicates from a previous conversation between this pair of hosts. The messages and packets must carry some identifying number for this purpose.

Conversations should be able to be broken by either and if an IMP finds its storage for message blocks filling up. The messages to do this, "do a reset" and "did a reset" as we call them, must also be error-controlled.

Conversations should begin without undue startup delay.

The dynamic tables should be simplex; that is, one-directional message numbers should be used to avoid deadlocks of the form A tries to talk to B, B tries to talk to C, C tries to talk to A, and none of A, B, or C have any more free blocks to use to begin a new conversation.

The tables should be fast to access at both the source and destination IMPs.

The method used the the ARPANET for implementing this system is based on a small pool of blocks, each of which carries a "use number", four bits wide, permanently associated with the block. All packets exchanged between IMPs carry the block number to be used in processing the packet and the use number. As explained in more detail below, these numbers provide the key information necessary for error control in a dynamic block environment.

The system works as follows: When a host at an IMP gives its IMP a message, the IMP first looks to see if a block exists for that source host to that destination IMP and host. Instead of seerdhing through all blocks, a "bucket sort" (a simple heah) is performed, to out the everage seerch length by some number (by it in actual fact) using a few bits of the key (source heat, destination IMP, destination host) to begin the sort. It then shecks successive blocks in the bucket led to by the sort (all the blocks in a single bucket are actually chained together for speed of access). If no block is found already existing for this key, a new block must be acquired at both the source and destination IMPs for this host/host conversation. The source IMP gets a block from the list of free blocks (the dase of no free blocks is discussed later). It puts the new block et the top of the bucket (head of

the chain) that it just smarched. The program then sopies in all the key information, adds one to the use number of the block, initializes the transmit message number entry, turns on a bit to indicate that the block is not yet in use, and calculates the index of the block it found.

The program than constructs a "get a block" massage which includes the index number calculated just above and the use number from the block and sends the "get a block" massage to the destination IMP. The sounce IMP then waits for an enswering "got a block". The "get a block" must be retransmitted every few seconds if no enswer returns. When a "got a block" is returned, the initialization bit is cleared, and the foreign block and use number which arrived in the "got a block" are copied into the source block. New the massage can be sent using the massage number window, allocation, and ordering techniques described above.

At the destination IMP, when a "get a block" is received, the program tries to get a free block. If the iree list is empty, nothing more is done (in effect, the "get a block" request is ignored). If a free block is aveilable, all the key data carried in the "get a block" is copied into the block and a "got a block" message is constructed including the key data and sent to the source.

When the source IMP sends a peaket to the destination, it carries the foreign block number and use number which are kept in the source block. The destination IMP uses this number to calculate the address of its message block and verifies the key information including the use number. In all other ways, the logic discussed above for accepting packets within the message number window is followed. When a Ready For Next Message is generated at the destination IMP, it services back the block number and use number kept in the destination block. This allows the source IMP to find its block and detect duplicate Ready for Next Messages in a simple manner.

We have explained how blocks are acquired. It is also necessary to discard blocks. There are two timers in the transmit and receive block. One counts two seconds of inactivity, the other two minutes of inactivity. A block may be discarded after two seconds of idle time, and a block must be discarded after two minutes of idle time, and a block must be discarded after two minutes of idle time. The two-second timer serves the function of quickly time-multiplexing the use of the dynamic blocks by many different conversations. The two-minute timer is used merely in a background menner to refresh the pool of free blocks by garbage sollecting blocks associated with conversations long inactive, thus avoiding more often expensive searches for a free block at the instant a new free block is required. The two-minute timer could be made longer if desired, to allow hosts to pause longer in

a conversation without incurring some setup delay; the two-second timer value is more oritical. If, as we assume, duplicate packets may arrive at an IMP up to thirty accords after initial arrival, then a mechanism is needed to allow blocks to be created and deleted more often then every thirty seconds that protects against the same pair of hosts using the same block twice and not catching a duplicate. The 4-bit use number allows sixteen cycles of acquisition and discard of the same block in any thirty second interval. Therefore, at least two seconds must elapse after acquisition of a block before it can be disparded. The rule that the message number must be (die for two segonds before a block den be discarded is actually stronger, but it seems a good rule to prevent threshing and inefficiencies. The two minute timer serves keep everything in synchronization in steady state (actually the timer runs for two minutes on the transmit side, and for signtly longer on the receive side, to prevent rages).

The best policy to follow for choosing when to delete blocks seems to be for an IMP to attempt to find deletable blocks (either transmit or receive) when its free block list goes below some slip, say 10% of the total peols. When this happens, if the program can locate a transmit block that can be deleted, it sends out a "reset" message to the destination, which causes the destination to discard its block and to send a "reset reply" back to the source which causes the source to discard its block. If

the program finds a requive block to delete, it sends a "reset request" message to the source which then follows the above protocol for performing a reset. On all these message, the block number and use number provide a duplicate detection facility, aince a given block with a particular use number can only be reset once.

At this point, it is worth noting the duplicate detection mechanism applied to the "get a block" and "got a block" messeges. The "get a block" carries no identifying information other than the addresses of the source and destination and the source bigsk number and use number. If a duplicate errives during the conversation it initiated, it can be detected: likewise, if it aprives during any other later conversation between those two hosts it can be detected. The only problem erises if the "get a block* duplicate aprives at the destination when no block exists between the two hosts. Then the destination IMP must get a block and return a "got a block" to the source. If the source receives a "got a block" when it is not expecting one, it must send out a Preset" to glear the destination. This fixes the problem, and Mince the program ignores Preset and Preset reply messages which do not match the blook and use number then ective, it also takes gere of the case of duplicate "got a block" messages and unwented *reset* and *reset reply* messages generated to des! With # + roumstandes.

This constudes our discussion of dynamic message blocks, They do not present a large penelty in storage, delay or packet they are reliable because they are maintained dynamically and all communications are error controlled; they allow the available blocks to be quickly multiplexed in time among a number of conversations (albeit at reduced performance for sonverset(on) rether then shutting some converset(one out as the available message transmission resources become fully used: they allow separate message numbers between each host/host pair even when the total number of hosts in the network is beyond a number where fixed linear tables indexed by host would be possible. In fact, expanding on the final point, the dynamic message block meshanism permits multiple message number streems between a given host/host pair. The ARPA implementation durrently permits two, one priority stream and one normal stream; but additional atreems are easy to imagine (e.g., special "permanent" streams which cannot be preempted after two seconds, streams which follow e different ordering or retrensmission griterion or even the lack of one, streams which have verying message number windows depending perhaps on stated host needs). Altogether, the dynamic message block mechanism appears to permit very flexible and reliable operation.

2.4.3.3.2 Restructured Message Numbers

The next major change involves not the packaging of message number tables into blocks but the restructuring of the actual message number tables themselves. The following changes are planned under this general heading:

a restructuring of RMESS, AMESS, and RALLY to simplify them and eliminate a redundancy!

expansion of the message window from 4 to 8;

extension of the priority bit concept to a general multi-level handling type, with independent message sequences for each type:

the depablity of always sending back the correct duplicate reply (RFNM, Dead, ste) in the event that the first reply is lost; this has not been implemented in the gurrent system; implementation of the message number scheme in such a way as to facilitate traffic through the network that does not use

The durrent program meintains RMESS, the next message number to accept, and AMESS, the next message number for which to send and answer (RFNM, ALLOCATE, etc), which are obviously not independent, with a window of 4, the rule is RMESS =< AMESS =< RMESS + 7, Further, a single bit is kept with RMESS to indicate if a message has been received, which gennot provide any

all (or any) of the measage-processing mechanisms.

information on exactly what state the message is in. Finally, the 4 bits in RALLY are used to specify the type of enswer to return. The basic idea, then, is to consolidate these data structures into some improved tables that serve the requirements better, Specifically, we can meet these goals with one message number and one set of status bits per message number.

The plan we have in mind is to keep RMESS as the next message number to accept (plus an offset, possibly), and use the bits in a new word, RSTATE, to indicate the atets of each message in the window, as follows:

idle: nothing deing; reply sent for this number = 8
request; got a request; allocate should be sent
message; eccepting this message
reply; got a message; reply should be sent

We will then use another word, RTYPE, to replace RALLY in detailing what kind of a reply to send (RFNM, RFN, VALL, DEAD, INC) in the idle and reply states and also to provide additional data of other kinds, so desired, in the request and message states, Figure 2-35 shows how RMESS and RSTATE interact, Note particularly the ambiguity of the RSTATE + I entry: it refers to RMESS + I unless it is a reply, then it refers to RMESS + I unless it is a reply, then it refers to RMESS + I unless it is a reply, then it refers to RMESS + I.

Figure 2=33 Restructured Message Numbers

Finally, the format of the message block gentaining the new restructured message numbers is shown in Figure 2-34,

Figure 2-34 Message Block Formets

ARPANET Completion Report
Chapter III Section 2 Design and Implementation

2.4 host-IMP

In the previous section we discussed a number of issues of end-toward procedure design which must be considered wherever the procedures are implemented, whether in the subnetwork or in the hosts. In this section we discuss the proper division of responsibility between the subnetwork and the hosts.

2.4.4.1 Basic Concepts

Extent of Message Processing in the packatment there has been considerable discussion in the packatment thing community about the amount and kind of message processing that should be done in communications subnetworks. An important part of the ARPA Network design which has become controversial is the ARPANET system of messages and packats within the subnetwork, ordering of messages, guaranteed message delivery, and so on. In particular, the idea has been put forth that such functions should reside at host level pather than subnetwork level.

We summerize the principle usually given for aliminating message processing from the communications subnetworks a) for complete reliability, hosts must do the same jobs, and therefore the nodes should not; b) host/host performance may be degraded by the nodes doing these jobs; a) network interconnections may be impeded by the nodes doing message processing; d) lockups can

happen in subnetwork message processing; a) the node would begome simpler and have more buffering gapacity (f it did not have to do message processing.

The last point is true, sithough the extent of simplification and the additional buffering is probably not significant, but we believe the other statements are subject to some question. We have previously given our detailed reasons for this belief. Here we simply summerize our main contentions about the place of message processing facilities in networks:

a. A layering of functions, a hierarchy of control, is essential in a complex network environment. For efficiency, nodes must control subnetwork resources, and hosts must dontrol host resources. For reliability, the basic subnetwork environment must be under the effective control of the node program == hosts should not be able to effect the usefulness of the network to other hosts. For maintainability, the fundamental message processing program should be nede software, which can be changed under central control and much more simply then all host programs. For debugging, a hierarchy of procedures is essential, single otherwise the solution of any network difficulty will require investigating all programs (including host programs) for possible involvement in the trouble,

b. The nature of the problem of message processing does not change if it is moved out of the natwork and into the hosts; the hosts would then have this very difficult job even if they do not went it.

- c. Moving this task into the hosts does not elleviate any network problems such as congestion, host interference, or suboptimal performance but, in fact, makes them worse since the hosts cannot control the use of node resources such as buffering, CPU bendwidth, and line bendwidth.
- d. It is besignily cheaper to do message prognosing in the nodes (smell (nexpensive computers) than in the hosts, and it has very few detrimental effects.

Recipiess. Reseasor. Connections. In a number of cases, an organization has desired to connect a large host to a network by inserting an additional minicomputer between the main host and the node. The general notion has been to locate the host-host transmission procedures in this additional machine, thus relieving the main host from coping with these tasks. Stated reasons for this notion includes

It is difficult to change the monitor in the main host, and new monitor releases by the host manufacturer poss continuing compatibility problems.

- . Core or timing iimitations exist in the main host.
- It is desirable to use I/O arrangements that may already exist or be available between the main host and the additional mini (and between the mini and the node) to avoid design or prosurement of new I/O gear for the main host.

While this approach may sound good in principle, and, in fact, may be the only possible approach in some instances, it often leads to problems.

First, the I/O arrangements between the main host and any preexisting peripheral processor were not designed for network connection and usually present timing and bandwidth constraints that greatly degrade performance. More seriously, the logical protocols that may have preexisted will almost certainly preclude the main host from acting as a general purpose host on the network. For instance, while initial requirements may only indicate a need for simple file transfers to a single distant point, requirements tend to change in the face of new facilities, and the network cannot then be used to full adventage.

Second, the peripheral processor and its software are often provided by an outside group, and the host organization may know even lose about their inner workings than they know about the main

host. The node is centrally meintained, improved, modified, and controlled by the Network Menager, but the peripheral processor, while an aqually foreign body, is not so fortunate. This issue alone is crucial; functions that do not belong in the main hosts belong in centrally monitored network equipment. Note that it is exactly those host groups who are unwilling to touch the main host's monitor who will be unlikely to be able to make subtle improvements in the protocols, error message handling and timing of the peripheral processor. From a broader esgnomic view, common functions belong in the network and should be designed once; the peripheral processor approach is a succession of costly special descendent the total goat is greatly escalated.

iong term solution to the dilemms is to have the verious The manufacturers support hardware and software interfaces connect to widely used networks. This is not likely to occur until commercial networks exist and are widely available. In the meentime, potential host organizations that wish to use early networks (like the ARPANET) should try to find ways to put connection directly into the main hest. enthropomorphic filustration may be helpfuls the network is, among other things, a set of standardized protogots or languages. potential network host is in the position of a person who needs to have dealings with people who speak a language he does not know, I) he does not want to learn the language, he can indeed choose to inconvenient, expensive, and unpleasant, and subtle meaning is always last. The situation is quite similar when a host tries to work through a peripheral processor. If a host wishes to interest with a network, it is usually unrealistic to try to make the host think that the network is a card reader or some other familiar peripheral. As usual, you get what you pay for.

Other Measure larvises. Other commonly suggested design requirement is for storage in the communications subnetwork, usually for messages which are currently undeliverable because a host or a line is down. This requirement should have no effect whatsoever on the design of the communications pert of the networks it is an orthogonal requirement which should be implemented by providing special storage hosts at strategic locations in the network. These can be at every mode, at a few modes, or at a single node, depending on the nelative importance of reliability, efficient line utilization, and cost.

Another commonly suggested design requirement is for the communications subnetwork to provide a message broadquat espablity; i.e., a host gives a message to its nade along with a list of host addresses and the nodes somehow send copies to all the hosts in the list. Again we believe that such a requirement should have no effect on the design of the communications part of

the network and that messages to be broadcast should be sent to a special host (perhaps one of the ones in the previous paragraph) for such broadcast.

2.4.4.4 Conclusions

We conclude this section on planned changes to the ARPA network with a brisf look at some modifications we would like to make to the interface between the IMP and the host. These plans are affected by the enormous difficulty of getting the host computers (of which there are more than 58) all to change their programs. Therefore, any changes we do make are likely to be installed so that they are compatible with existing conventions and formets. By this we mean that the IMP and host will have to identify which conventions are in use, and the IMP will have to keep the apitwere to support both formets for an erbitrary length of time.

One simple kind of change, yet a difficult one to implement everywhere, is an extension of many of the fields in the message leader, to allow more than 64 IMPs and more than four hosts per IMP. We would also like to add other fields to the leader, or extend the meaning of existing fields. For example, the idea that the host interface should not be stopped without the host's consent means new message types must be defined.

He also feel that the host interfaces should be made to run as fast as possible. This is not so much to allow hosts to get higher throughput, since they will still be limited by network sirguit bandwidths (though it will speed up intracIMP messages), but to out down delays on host lines. This would make it easier for us to make enother changes not sending back the RFNM until the whole message is accepted by the host. Currently, the RFNM is generated after the host accepts the first packet. This is done for two reasons. First, it is assumed that the host-IMP interface is error-free and further that once the host begins to accept bits of the message it will accept it all. Secondly, to delay the RFNM until the and of the message would add tens or hundreds of milliseconds of delay to the roundetrip time. Speeding up the host interfaces would reduce this delay, and we would prefer to make the RFNM mean that all the message was delivered.

Again on the subject of reliability and addurate message ecounting, we have mentioned that endatowend checksumming would also be desirable. Perhaps error detection and netransmission over the IMP-host interface is the right way, perhaps the IMP should denerate and verify checksums, perhaps the hosts should. This is really a subject for negotiation among the various parties. It is our basis belief that the subnetwork at least should provide such a message-processing facility, and that hosts an augment it if they desire.

We have also some around to believing that the information that the IMP and host can exchange about the status of specific messages and their states generally, the better the network is to use. Therefore, the IMP now reports to the host. ss closely as (t can why a message sould not be delivered,, distinguishing about a dozen cases at present, Further, we have implemented thost going donw! and 'IMP going down' messages, which carry information on why the system is going down, when it is scheduled to go down, and when it is scheduled to be back up. This information can also be applied on a message-by-message basis, to supplement any non-delivery diagnostics. Generally, feel such information machenisms are an important adjunct to the banic message-processing facilities. They greatly facilitate the meintanence, and debugging of the system associated with the ARPA network and we expect that we may see more of this kind of diagnostic information passed around at all levels of the networks

2.5 Howt Lavel Protocols

The ARPANET subnet of IMPs and communication paths, as described in the previous sections, provides the physical mechanism which allows the transport of messages from one location to another. The IMP/Heet interface, both hardware and software, provides the mechanism which allows the submission of a message to the network by one Host and the reception of the message from the network by another Host. Nevertheless, this set of abilities is insufficient to insure the attainment of meeningful communication among domputer programs in the various Hosts.

An illustrative analogy may be made with the telephone system; the existence of the switching system, connecting handsets in (simpst) every home and office, is insufficient to insure meaningful communication between an English speaker and a French speaker, or even among several English speakers connected in a "conference cell". There must be rules of discourse to insure both a common choice of language, and a common mode of interaction, to make communication meaningful.

In the ARPANET the rules of discourse which govern the exchange of information between two (generally) similar entities

are denoted protosols. Protosols deal not only with the details of information format, the dommon language, but also with the framework (or model) within which information is exchanged. Of course, any pair (or larger group) of cooperating entities may agree to a particular set of rules of discourse, but the protocols of interest are those which were adopted by substantial segments of the ARPANET dommunity. This is not so because the "standard" protocols are necessarily better, but because implementation of a standard protocol admits a Host to a large dommunity of communicating entities, while implementation of pairwise non-standard protocols allows communication only with other entities which follow the same non-standard protocol.

Naturally, there have been situations within the ARPANET community where development of an admind protogol was the most advantageous course for a limited user group. These situations commonly arise from a desire for increased efficiency along some dimension (e.g., implementation size or transmission sped) and the solutions chosen frequently obtain an efficiency by giving up some measure of generality provided by the standard protocol. A few examples of such adehog protocols which have been implemented within the ARPANET environment includes

- COPYNET, an interim protocol for file transfer between PDP=10 TENEX systems in 36=bit word format. The protocol was discontinued with the adoption of a standard file transfer protocol.
- A file transfer protocol adapted to achieve high bendwidth between two Univer 418=III computers at Tinker and McCiellan Air Force Bases. The existing Univer record structures were transformed directly into ARPANET messages.
- A file transfer protocol adapted to achieve low memory requirements for TIP implementation; implemented on TIPs equipped with the magnetic tape option.
- A real-time message transfer protocol used between setamic data collection sites and a centralized setsmic data processor.
- A resisting message transfer protocol used for the experimental transmission of digitized speech.

Nevertheless, the vest mejority of traffic which has been carried by the ARPANET has been handled by the standard protocols.

ARPANET Host level protogol dealgn has followed a The of layering, a strategy which is similar to the etratedy levels-of-ebstraction approach to program development. laver, or level, of Host protocol builds upon the lower lavers. This strategy has several important potential advantages. First, such layer of protocol is assigned a clearwout set of functions; designers of each laver are presented with a bounded set of objectives and are able to think about solutions within this bounded context. It is reasonable to assume that this makes the design process easier, and hence leads to better overall design. Second, each lever of protocol can present a straight-forward functional interface to the next higher layer. This means that each layer may be tested independent of higher layers, which should lead to easier and quicker implementation. It should also mean that a functional sepablifty contained within one layer can be redesigned and/or resimplemented without somplications arising in other lavers.

On the other hand, there is a danger that too many layers may be created; if data must gross many interfaces, each requiring some transformation, on its path from the original producer to the ultimate consumer than the system as a whole may be extremely inefficient. There have been some examples in the

ARPANET experience of inefficiency resulting from the layered design; for an example see Section 2.5.1 which briefly discusses the official RJE protocol.

The layering structure adopted for the ARPANET Host level protocols is based upon the model of the way in which a typical time-shared computer might be structured. At the lowest level is a monitor to monitor (or operating system to operating system) protocol, which is denoted the Hoste-Host protocol. This layer is envisioned as handling communication issues such as multiplexing many independent communications over a single interface, or managing I/O buffers.

The next layer is denoted the Initial Connection Protocol and addresses the issue of getting a Host's monitor (operating system) to enter a new user (nto the current user set. This is, obviously, a transitory issue; the "Initial Connection Protocol" is similarly a transitory layer, it is a procedural protocol mather than a formatting protocol.

The third layer of Host level protogol deals with the method of terminal-to-Host (or user to application program) dommunication. The protogol is denoted "TELNET", a word derived from the phrase "Telecommunications Network". Originally it was

envisioned that this protosol would be one of many "user level protocols" in the same lever, but as the system model of the structure of a time shering system gained strength it began to seem appropriate to use this protocol as the "dommand path" building block for other protocols, and thus TELNET became a layer in its own right.

The fourth layer of protocol is the File Transfer Protocol, This is, as indicated by the name, the specification of a mechanism for the transfer of complete files, not a mechanism for retrieval of some information subset. The File Transfer Protocol (FTP) uses the TELNET protocol to govern the exchanges of commands and responses which specify the transfer, and uses the Mosta-Host protocol directly for the actual movement of the file.

At a still higher layer is the ARPANET official Remote Job Entry (RJE) protocol, which is designed to handle the submission of batch jobs and retrieval of results. This protocol uses TELNET to govern the exchange of Job control commands and responses, and uses FTP to govern the actual transfer of Job inputs and outputs. In actuality, there are probably few (or no) implementations of the Mofficial RJE protocol, most ARPANET batch Job submissions use the ad hog RJS (Remote Job Service)

protocol specified by UCLA for use with the IBM 360/91 system located there.

Figure 2.5=1 below (liustrates the layering of protoco) graphically. The following subsections briefly describe each of the Host level layers in greater detail.

Figure 2:5-1: Protocol Layering

2.5.1 Hostwhost Protocol

The ARPANET Hostodost protocol is the lowest layer of the Host level protocols. As noted above, it is envisioned as addressing those issues relevant to the Host system taken as a whole, that is, the issues of operating system to operating system communication.

2.5.1.1 Basia Concepts

The Host-Host protocol for the ARPANET was designed to address the following five basic (sauce:

- Multiplexing: It was envisioned that there would be many processes within a single Host all trying to communicate through the network simultaneously. The Host-Host protocol was required to provide a way of multiplexing these communications onto the single Host-to-network interface.
- ErrailCantral: The Host-Host protocol designers had to choose the appropriate level of seror control, both in terms of bit errors, and also in terms of lost, duplicated, or disordered messages.

- Eigh__Casteli The Host-Host protect should include medhanisms which would prevent a fast transmitter from overrunning the processing or buffering capacity of a slow receiver.
- Quitarizand.Signalling: It was anticipated that in any communication system with flow control there may be president when critical information is prohibited from being transmitted due to the flow control rules. In such cases an "out-of-band" (i.e., not subject to flow control) signalling system may be desirable, if not essential.
- Respectable in communicating systems may occasionally get out of steps with the "state information" describing the surrent situation different in the two systems due to herdware or software failure somewhere. The Host-Hest protocol should provide a way of resynchronizing the two systems.

In addition to addressing these five basic issues, which can be seen recurring in a wide veriety of protocol design situations, the ARPANET protocol designers were influenced by several other goals and consepts. First, the protocol was

designed specifically for the ARPANET environment. In practical terms, the concentration on this environment had three significant effects on the protocol design, as follows:

- The protocol decidners accepted at face value the claims of the subnetwork implementars that the network would be error-free, Because of this, no bit or message error detection or recovery mechanisms were included in the protocol.
- The IMP/Host interface specification left a few bits in the message leader (i.e., the fixed format portion of the message in which the message address is included) to be assigned by the sending Host; this was originally delied the "link" field end is now a portion of what is delied the "message ID" field. The Host=Host protocol design attempted to insure that no information which wouldn't fit into this field would be required on a per-message basis, in an effort to avoid adding overhead.
- 3) The IMP subnetwork was apacified to deliver a special control message (the Reedy for Next Message; or RPNM) to a message sender to serve both for flow sontrol at the

Host-to-IMP level and to acknowledge that the message had been successfully transported across the subnatwork. The Host chose to treat this as a surrogate endeto-end acknowledgement.

Each of these instances of reliance upon the (actual or assumed) characteristics of the IMP subnetwork had a positive effect on the ease with which the protocol could be initially specified or implemented, but has had some longerance negative effects, as will be discussed below.

A second concept which had a strong influence on the design was the model chosen as a framework grotocol dommunication. The model chosen was the familiar system of sircuit switched telephone communication, in an implementation of what have some to be salled "virtual singuite" in other contextas the entities are called "connections" in the ARPANET context, Thus, the ARPANET subnetwork, which "switches" individual messages, is used by the Host-Host protosols as a base upon which to establish dirduits. (It is interesting to note that earlier in the history of computer communications significant effort has been invested in developing techniques for building message switching systems on a base consisting of direults: one obvious

example is the development of "multipoint" lessed line management.) An alternative choice of model (e.g., the postal system) might have led to a completely different protocol. In fact, some experiments with a message-based Hoste-Host protocol were suggested fairly early in the ARPANET development, but none were actually implemented.

Third, given a model of sommunication which was based on "connections", the protocol design was strongly influenced by a desire for dynamic reconnection, or handwolf of one and of a connection from one processes to another, with the previous termination point and the new termination point (potentially) in different Hosts. As described later, this goal led to the necessity for the complex Initial Connection Protocol and influenced the specification of simplex, rether than full-sduplex, connections. In retrospect, this seems to have been a high price to pay, since it appears that the potential for dynamic reconnection has never been used by any Host.

Another fairly basic idea which grew out of the connection-based model was the concept of a "socket". A socket is a particular implementation of a more general problem; the creation of a name apace for processes which is global in scope,

It was obvious that each individual Host would have some internal scheme for naming processes, but these schemes were enticipated to be at least different, and possibly even incompatible. Since it would be impracticed to try to impose a common internal process naming scheme, an external name space was greated, with a portion essigned to each Host; this name space consisted of sockets, which were specified by condetenating a Host address with a 32-bit "socket number". Obviously, the same socket number at two different Hosts represented two different processes. Each Host was given the responsibility for mapping internal process identifiers into logal socket numbers.

It was envisioned by the protocol designers that there would be at least two methods of assigning socket numbers to processes some socket numbers would be assigned permanently to processes which requested them. For example, a service provided by the Host would need at least one socket which remained gonstant so that potential users would know how to address requests for service. A programmer might be assigned a group of sockets which sould be used by his programs; the programmer would be able to convey these numbers to other programmers at other sites who might want to communicate with his programs. On the other hand, some socket numbers might be assigned by an operating system on a

transitory basis to processes which wanted to send service requests into the network. A socket assigned in this way gould be relinquished explicitly when the process had no further summunication needs, or eutomatically when the process terminated. In fact, both methods of socket assignment are common in the ARPANET.

A connection is specified by naming the two sockets which it connects. However, carrying 64 bits of socket number in each message was deemed by the protogol designers to be an exactsive amount of overhead. In order to reduce the amount of overhead, it was decided that use of a connection should be preseded by a setup phase (in fact, it is primarily this setup phase which differentiates a connection from a stream of messages) during which an "abbreviation" for the 64 bits of connection name would be established. As noted above, the IMP/Host interface specified message leader which contained an 8-bit field. (originally) gailed the "link", which was to be specified by the sending. Host and which would be delivered by the IMP aubnetwork elong with the measage text. It was elder that if a mapping from 54 bit gonnection identification to 8=bit link were constrained to be unique at any instant, the link could serve as an abbreviation for the connection identification, and the connection

DRAFT

identification would then require no new message overhead. Thus, protesol is closely wedded to the concept of an the HosteHast IMP/Host link.

Unfortunately, the concept of the IMP/Host 1 ink originally created for purposes of flow control (sithough, as has been noted in a previous section, it was unsultable for this purpose and eventually abandoned by the IMPs). The IMPs therefore imposed the rule that a Host dould not have two messages with the same link number in the network at the same time. This being the gase, the RFNM (or less frequently, the Incomplete Transmission response, which identified messede) could refer to a link number and uniquely identify a message. Thus, use of the link to identify a connection, and relying on the RFNM/Incomplete as the only messegement and accorcontrol, meant that a HosteHost connection could never have more than one message in flight at a time, even after the IMP subnetwork removed this restriction. This has been an extremely important impediment to achieving good bandwidth performance between Hosts. Waiting for a RFNM between messages typically limits connection bendwidth to ebout 18Kbs, in contrast to the 40 Kbs or so which is usually available from the subnetwork,

Another consequence of the decision to use the IMP/Host link field in this way was the decision to make connections simplex, rether than full-duplex. Since the link field was used to identify the connection to which a message belonged, it had to be unique. Since it would be impossible to guarentee that when two Hosts established a connection they would be able to find any link number which was free at both ends, it was clear that each Host should choose the link number which it would use to identify the conversation. However, since the IMPs required that not more than a few tens of links be in use at one time, link numbers were viewed as a scerce resource, and thus silocating one to an unnecessary date path seemed wasteful. For example, a file transfer would typically need only a simplex connection,

Of course, if each connection must go through a setup phase before it exists, there must be a mechanism in each Host which dan manage the setup phase and, further, these mechanisms must be able to communicate before the first connection is set up. The ARPANET Host-Host protocol envisions a process delied the Network Control Program (NCP) running in each Host; this process menages the IMP/Host interface (this is the sense in which it "controls" the network) and also manages the Host-Host protocol. It is the NCP which performs connection setup, Messages sent from one NCP

to enother for control purposes are referred to an "control messages" and are identified by having the link number set to sero; thus, link sero is denoted the "control link," It is useful to think of the control link as implementing a control sennection which is a permanent, rather than a "switched", connection.

2.5.1.2 Original Specification

The Host-Host protogol is specified in detail by the specification of the effects of the control commends exchanged between NCPs. The specific functions of the NCPs which are effected by these commends are the establishment and termination of connections, the control of the rate of date flow across established connections, and the transmission of outerfaces lit should be noted that an NCP has at least two interfaces; one to the external world, which is specified by the Host-Host protocol, and a second to the processes, running in the Host, which the NCP has serving. This internal interface was considered a purely Host-specific issue, and the standard protocol made no ottompt to standardize it. Nevertheless, it is possible to describe typical instances of this interface in general terms; such a description is given in Section 2.5.1.2.6.

2.5.1.2.1 Connection Establishment

Two commands are defined to establish a connection; these are STR (sender-to-receiver) and RTS (receiver-to-sender), each of which takes three parameters:

STR (send socket, receive socket, byte size)

RTS (receive socket, send socket, link)

The STR command is sent from a prospective sender to a prospective receiver, and the RTS from a prospective receiver to a prospective sender. The send sopket field names a socket local to the prospective sender; the requive socket field names a socket local to the prospective requiver. In the STR command, the "size" field contains an unsigned binary number (in the range to 25%; zero is prohibited) specifying the byte size to be used for all messages over the connection. In the RTS command, the "link" field specifies a link number; all messages over the connection must be sent over the link specified by this number; These two commands are referred to as requests for connection (RFCs). An STR and an RTS metch if the receive socket fields match and the send socket fields match. A connection is established when a matching pair of RFCs heve been exchanged.

With respect to a perticular connection, the Host containing the end socket is called the "sending Host" and the Host containing the require socket is called the "requiring Host". A Host may connect one of its receive sockets to one of its send sockets, thus becoming both the sending Host and the receiving Host for that connection. These terms apply only to data flow; control messages will, in gameral, be trensmitted in both directions.

Hosts are prohibited from establishing more than one connection to any local socket. This restriction has led to rether severe difficulties in establishing contact with a general service, as exemplified by the Initial Connection Protocol (see Section 2.5.2). It has been argued that in general only one socket of the pair which specify a connection need be unique in order to evoid confusion; the stronger restriction stated for the ARPANET Hostshost protocol is a result of the desire for a dynamic reconnection mechanism which "slways" works. As an example, consider the case in which sockets a and y in Host A were both connected to socket a in Host B. This situation causes no ambiguity. However, if all of the connections emanating from Host A needed to be temporarily reconnected to some new location, say socket with Host C, then there would be two connections from

ORAFT

w (at C) to z (at 8); it would no longer be possible to distinguish them. The problem arises whenever it is necessary to reconnect connections terminating at two general service process. each of which "assumes" that the Host at the other end will take same to keep its and unique. To avoid over making such an assumption, the atronger condition is imposed.

A Host sends an RFC either to request a connection, or to accept a foreign Hoat's request. Since RFC commends are used both for requesting and for accepting the establishment of a connection, it is possible for either of two cooperating initiate connection establishment. to consequence, a family of processes may be drested sonnection-initiating actions built-in, and the processes within this family may be started up (in different Hosts) in arbitrary order provided that appropriate queueing is performed by the Heats involved.

There is no prescribed lifetime for an RFC. A Host is permitted to queue incoming RFCs and withhold a response for an arbitrarily long time, or, alternatively, to reject requests (see Connection Termination below) immediately if it does not have a matching RFC outstanding. It may be reasonable, for example, for an NCP to queue an RFC that refers to some currently unused agaket Until a local process takes control of that socket number and tails the NCP to accept or reject the request. Of course, the Host which sent the RFC may be unwilling to wait for an arbitrarily long time, so it may abort the request. On the other hand, some NCP implementations may not include any space for queueing RFCs, and thus can be expected to reject RFCs unless the RFC sequence was initiated locally.

The decision to queue, or not queue, incoming RFCs has important implications. Each RFC which is queued, of gourse, requires a small amount of memory in the Host doing the queueing. If each incoming RFC is queued until a local process seizes the local socket and accepts (or rejects) the RFC, but no local process ever seizes the socket, the RFC must be queued "forever," Theoretically this could occur infinitely many times (there is no reason not to queue several RFCs for a single local socket, letting the local process decide which, if any, to accept) thus requiring infinite storage for the RFC queue, On the other hand, if no queueing is performed sooperating processes will be able to establish a desired connection only by accident (when they are started up such that one issues its RFC while the RFC of the other is in transit in the network = glearly an unlikely occurrence).

Perhaps the most ressonable solution to these problems is for each NCP to give processes running in its own Host two options for attempting to initiate connections. The first option would allow a process to cause an RFC to be sent to a specified remote socket; with the NCP notifying the process on to whether the RFC were assented or rejected by the remote Host. The second option would allow a process to tell its own NCP to "listen" for an RFC to a specified local socket from some remote socket (the process might also specify the particular remote socket and/or Host it wishes to communicate with) and to secept the RFC (i.e., return a matching RFC) if and when it arrives. Note that this also involves queueing (of "listen" requests), but it is internal queueing which is susceptible to reasonable management by the local Host, If this implementation were available, one of two cooperating processes could "listen" while the other caused a series of RFCs to be sent to the "listening" socket until one was accepted. Thus, no queueing of incoming RFCs would be required, although it would do no harm.

2,5,1,2,2 Connection Termination

The command used to terminate a connection is CLS (closs), which takes two parameters:

CLS (my socket, your socket)

The "my socket" field contains the socket local to the sender of the CLS command. The "your socket" field contains the socket local to the receiver of the CLS command. Each side must send and receive a CLS command before connection termination is completed and the sockets are free to participate in other connections.

It is not necessary for a connection to be established (i.e., for both RFCs to be exchanged) before connection termination begins, For example, if a Host wishes to refuse a request for connection, it sends back a CLS instead of a matching RFC. The refusing Host then waits for the initiating Host to scknowledge the refusal by returning a CLS, Similarly, if a Host wishes to abort its outstanding request for a connection, it sends a CLS command, The foreign Host is obliged to acknowledge the CLS with its own CLS; even though the connection was never established, CLS commands must be exchanged before the sockets are free for other use,

After a connection is established, CLS commands sent by the receiver and sender have slightly different effects. CLS sommands sent by the sender indicate that no more messages will

be sent over the connection; thus the CLS must not be sent if there is a message in transit over the connection. A CLS command sent by the receiver acts as a demand on the conder to terminate transmission. However, since there is a delay in getting the CLS command to the sender, the receiver must expect more input.

A Host should "quigkly" anknowledge an incoming GLS so the foreign Host can purge its tables. However, there is no prescribed time period in which a GLS must be acknowledged.

Because the CLS command is used both to initiate closing, aborting and refusing a connection, and to agknowledge closing, aborting and refusing a connection, race conditions can occur, However, they do not lead to ambiguous or arroneous results, as illustrated in the following examples:

EXAMPLE 11 Suppose that Host A sends B a request for connection, and then A sends a CLS to Host B because it is tired of waiting for a reply. However, Just when A sends its CLS to B, B sends a CLS to A to refuse the connection. A will "believe" B is acknowledging the abort, and B will "believe" A is acknowledging its refuse; but the outcome will be correct.

EXAMPLE 2: Suppose that Host A sends Host B an RFC followed by a CLS as in exemple 1. In this case, however, B sends a metching RFC to A just when A sends its CLS, Host A may "believe" that the RFC is an attempt (on the part of B) to establish a new connection or may understand the race condition; in either case it can discard the RFC since its socket is not yet free, Host B will "believe" that the CLS is breaking an established connection, but the outcome is correct since a matching CLS is the required response, and both A end B will then terminate the connection.

Every NCP implementation is feed with the problem of what to do if a metching CLS is not returned "quickly" by a foreign Host (i.e., if the foreign Host appears to be violating protocol in this respect). One naive answer is to hold the connection in a partially closed state "forever" waiting for a metching CLS. There are two difficulties with this solution. First, the socket involved may be a "searce resource" such as the "logger" socket specified by an Initial Connection Protocol (see Section 2,5,2) which the local Host cannot afford to tie up indefinitely. Second, a partially broken (or mailing out) process in a foreign Host may send an unending stream of RFCs which the local Host wishes to refuse.

Most implementars, recognizing these problems, have adopted some unofficial timeout period after which they "forget" a connection even if a metahing CLS has not been requived. It has been suggested that the network adopt some standard timeout period, but it has been proven impossible to arrive at a period which is both short enough to be useful and long enough to be acceptable to every Host. The difficulty with any timeout is that if a second connection between the same pair of sockets is later established, and a CLS finally arrives for the first connection, the second connection is likely to be closed. This situation can only arise, however, if one Host violates protocol in two ways; first by failing to respond quickly to an incoming CLS, and second by permitting establishment of a connection involving a socket which it believes is already in use.

2.5.1.2.3 Flow Control

After a connection is established, a sending Host sends messages over the agreed-upon link to the receiving Host. The receiving NCP accepts messages from its IMP and queues them for its various processes. Since it may happen that the messages arrive faster than they can be processed, some mechanism is required which permits the receiving Host to quench the flows from the sending Hosts.

The flow control mechanism requires the receiving Host to allocate buffer space for each connection and to notify the sending Host of how much space is available. The sending Host keeps track of how much room is available and never sends more data than it believes the receiving Host can accept.

To implement this mechanism, the sending Host keeps a bit downter essociated with each connection. The counter is initialized to zero when the donnection is established and is increased by allocate (ALL) control dommands sent from the receiving Host as described below. When sending a message, the NCP of the sending Host subtreats the "text length", i.e. the product of the connection byte size and the number of data bytes in the message, from the bit counter. The sender is prohibited from sending if the counter would be degreemented below zero.

The flow control mechanism does not pertain to the control link, since connections are never explicitly established over this link.

The control command used to increase the mender's bit sounter was ALL (allocate), which took two persmeters:

ALL (link, bit space)

This command is sent only from the receiving Host to the sending Host, and is legal only when a connection using the link number appearing in the "link" field is established. The "bit space" field is defined to be an unsigned binery integer specifying the amount by which the sender's bit counter is to be incremented. The receiver is prohibited from incrementing the sender's counter above 2 to the 32nd minus i. In general, this rule will require the receiver to maintain a counter which is incremented and degreenested according to the seme rules as the sender's counter.

2.5,1.2.4 Out=of=Band Signalling (Interrupt)

The Host-Host protocol includes a meshenism by which the transmission over a connection may be "interrupted." The meaning of the "interrupt" is not defined at this level, but is made eveilable for use by outer layers of protocols. The interrupt command sent from the neceiving Host to the sending Host is INR (interrupt-by-receiver).

INR (1(nk)

The interrupt sommend sent from the sending Host to the receiving Host is INS (interrupt=by=sender).

INS (1+nk)

The INR and INS commands are legal only when a connection using the link number in the "link" field is established.

2.5.1.2.5 Other functions

In spite of the effort made to insure that the HosteHost protocol functions did not require any edditional bits to be transmitted on a peramessage basis (using instead the IMP-Host protocol "link" field), it seemed slear that many Hosts might have difficulty beginning a message immediately following the 32-bit IMP/Host leader. For example, a PDP=10 Host, with a 36-bit word size, would find it much more convenient to begin the ectual text of a message at a 36mb/t boundary; beginning at a 32-bit boundary would be likely to require extensive expensive bit shifting. Therefore, the protocol defined a "marking" field to be included in each message between the IMP/Host leader and the first meaningful data bit. The marking consists of zero or more zero bits followed by a single one bit. Thus, the receiving Host NCP begins at bit 33 of a received message searching for a bit. The bit following this one bit is the beginning of the meaningful deta.

It is sametimes useful for one Host to determine if some other Host is depeble of carrying on network conversations. The control command to be used for this purpose is ECO (scho).

ECO (data)

The "data" field may contain any bit configuration chosen by the Host sending the ECO. Upon receiving an ECO command an NCP must respond by returning the data to the sender in an ERP (Egho-reply) command.

ERP (data)

A Host should "quickly" respond (with an ERP command) to an incoming ECO command. However, there is no prescribed time period, efter the receipt of an ECO, in which the ERP must be returned. A Host is prohibited from sending an ERP when no ECO has been received, or from sending an ECO to a Host while a previous ECO to that Host remains unenswered.

2.5.1.2.6 NCP/Process Interface

As previously noted, the interface between an NCP and the processes which use it within the same Host was considered to be an internal matter, not a subject for network-wide

standardization. Nevertheless, the design of the Host-Host protocol assumed that certain facilities would be available to the processes which used the NCP, and these were documented in terms of a set of hypothetical system dails which NCP implementers might provide.

atrongly influenced by the facilities supported by the Host operating system. An operating system consists of program modules which augment the hardware and provide an environment for processes. The interfaces between these operating system modules and user processes take the form of system calls and returns, and sometimes pseudo interrupts. System calls are implemented in a variety of ways, but often it is a special hardware instruction that involves a system call (e.g., SVC, UUC, JSYS, MME). In higher level programming languages a system call is often indistinguishable from a subroutine call. In some cases the form of the system call is different for each module of the system.

A product is a program in execution with its associated address space, a location counter, some general registers, and usually some open files (or devices). Prodesses may be created by users, though there are often processes which have been

areated as part of the eystem for particular functions, and some of these may be initialized by the system when it begins running. A process may have several entry points, each of these is called a Port. Special Ports which are designed in processes to handle unusual and error control conditions are called "code" returns.

One important expect of an operating system that has a great impact on the implementation of network functional capabilities is the provision for interprocess communication. Generally processes are viewed as independent computational units that need interest only with the operating system in a few very constrained ways (i.e., via system cells). However, often it would be useful to build a new capability based on a dombination of existing programs. One way of extending the usefulness of the process structure is to allow processes to communicate among themselves such that several efstagods vem assassong to eccomplish a computational goal. The form of communication aupported by an operating system very much influences the extent to which processes actually cooperate and, therefore, the extent to which ARPANET is a natural extension of the programming use of the environment.

A basic set of system calls which should be available to user written processes for use of the natwork via the Hastwhost protocol are:

LISTEN (port, socket, ende): This is a cell to the NCP indicating that the process is willing to connect to any process which sends a request for connection (RFC) specifying this local socket. This requests the NCP to essociate a local socket with a process and the specified port in that process. A return value is given in "code", If there is a pending gall, the connection may be opened immediately and the process notified; if there is no matching pending call, the NCP will notify the process when a matching RPC errives.

CONNECT (port, modket, formign mocket, dode); The local socket of this process is especiated with this call. The local process requests the NCP to initiate a connection on its behalf to a particular foreign mocket from the specified local mocket and mesodiate the donnection with the specified port. If there is a pending dail matching these parameters, the donnection is opened and the process so notified. A return value is given in "code". If there is no pending

matching call, the NCP communicates this request to the foreign host and notifies the local process when a matching request is received and the connection is opened.

- SEND (port, buffer, length, gode); The logal process requests the NCP to transmit the data starting at "buffer" and extending "length" bits over the connection associated with this port in accordance with the allocation values, "Code" is set with a return value.
- RECEIVE (port, buffer, length, code); A local process requests the NCP to store data received on the gonnection associated with this part into the processes address space starting at "buffer" and extending for "length" bits. A return value is set in "gode".
- CLOSE (port, gode): Activity on the connection essociated with this port is stopped. A return value is set in "code".
- INTERRUPT (port, gode): The local process requests the NGP to send a special interrupt signal referring to the connection associated with this port outside the normal flow control constraints. A return value is set in "code".

STATUS (port, info, code); This is a call by a local process requesting the NCP to obtain the relevant status information from the connection table entry associated with the port and place it in the space specified by "info", A return value is set in "code". This allows a user program to monitor the state of a connection, of special interest are the flow control allocation values.

2,5,1,3 Subsequent Modifications

Within a few months of the specification of the Host-Host protocol, as the first few implementations began to take form, several shortcomings or inefficiencies deme to light. A revised version of the protocol was developed to dorrect for these shortcomings; it is this revised version which has survived as the standard Host-Host protocol. The modifications made were in the areas of flow control, NCP synchronization, and text alignment (*marking*); each of which is described briefly.

Riow_sested_shappes. The early implementers identified two types of difficulty with the flow control mechanism. First, some Hosts with limited sapphilities wished to allocate only a single buffer to each natwork connection, and to avoid the bit=packing which would be needed if multiple short messages (whose total

length was less than the number of bits in the buffer) arrived, The solution adopted was to add a message count to the ALL commend. This allows a Host to restrict the total number of bits it is willing to receive, the total number of messages, or both. Thus, the form of the ALL commend is now

ALL (link, message spage, bit spage)

A second criticism of the flow control scheme is that it requires a Host to explicitly divide its total buffer space among all the connections which it may possibly participate in. It was argued that it would be better if a Host could assign its total buffer apace to the set of existing connections (or even "overbook" its ellocations), retracting some apace sllocation for reassignment if a new connection is established or, in the case of overbooking, if the space is being used up more rapidly then expected. To accomplish the retraction of allocation, two new commands were defined. The control command used to request that the sending Host return all or part of its current allocation is GVB (give=back), which takes three parameters:

GVB (link, message fraction, bit fraction)

This command is sent only from the receiving Host to the sending Host, and is legal only when a connection using the link number in the "link" field is established. The message fraction and bit frection fields are defined as the fraction (in 128ths) of the surrent message space allocation and bit space allocation (respectively) to be returned. If either of the fractions is equal to or greater than one, all of the corresponding allocation must be returned. Fractions are used since, with messages in transit, the sender and receiver may not agree on the actual ellocation at every point in time.

Upon receiving a GV8 command, the sending Host must return at least the requested portions of the message and bit space allocations. A sending Host is prohibited from spontaneously returning partions of the allocations. The gantrol command for performing this function is RET (return), which takes three parameters:

RET (|ink, message space, bit space)

This command is sent only from the sending Host to the receiving Host, and is legal only when a connection using the link number in the "link" field is established and a GVB command has been received from the receiving Host. The "message space" field and

the "bit space" field are defined as unsigned binary (ntegers specifying the amounts by which the sender's message counter and bit dounter (respectively) have been degreemented due to the RET activity (i.e., the amounts of message and bit space allocation being returned). NCPs are obliged to answer a GVB with a RET "quickly"; however, there is no prescribed time period in which the answering RET must be sent.

It is interesting to note that no Hosts assually adopted a buffer essignment strategy which requires the use of the GVB.

MCS. System "creshes", NCP errors, the addition of new network Hosts, or other feators, communication between two NCPs will need to be synchronized. One possible effect of any disruption might be that neither of the involved NCPs could be sure that its stored information regarding connections with the other Host matched the information stored by the NCP of the other Host. In this situation, an NCP may wish to reinitialize its tables and request that the other Host do likewise; for this purpose the pair of control commands RST (reset) and RRP (reset=reply), neither of which take parameters, were added to the protocol.

The R&T command is to be interpreted by the Host receiving it as a signal to purpe its NCP tables of any entries which arose from communication with the Host which sent the R&T. The Host sending the R&T should likewise purpe its NCP tables of any entries which erise from communication with the Host to which the R&T was sent. The Host receiving the R&T should acknowledge receipt by returning an RRP. Once the first Host has sent an R&T to the second Host, the first Host is not obliged to communicate with the second Host, the first Host is not obliged to communicate with the second Host (except for responding to R&T) until the second Host returns an RRP. In feat, to evoid synchronization errors, the first Host should not communicate with the second until the R&T is answered. If both NCPs decide to send R&Ts at approximately the same time, then each Host will receive an R&T and each must answered.

Hosts ere prohibited from sending an RRP when no RST has been received. Further, Hosts may send only one RST in a single sontrol message and should wait a "reasonable time" before sending another RST to the same Host. Under these conditions, a single RRP constitutes an "answer" to all RSTs sent to that Host, and eny other RRPs arriving from that Host should be discarded.

Text. Alignment: It was quickly regognized that the marking sonvention adopted for text alignment completely solved the sender * problems, but didn't help the receiver at all. In view of the fast that all of the Hosta connected (or anticipated to be sonnected) to the ARPANET sould conveniently handle units of 8. 18, or 36 bits, it was decided that all messages should begin with a "message header" of 72 bits (the least common multiple of these lengths) including the IMP leader. Singe this adds 40 bits to the length of every message, it was desided to specify that 8 bits are used to restate the number of bits per byte (as established during connection setup) and 16 bits are used to specify the number of bytes in the message. These fields ere used to locate the last bit of the text,

The protocol es originally specified, with the modifications described ebove, has been the "official" atandard Hostwhost protocol for the ARPANET einse early 1971. Because of the large number of different implementations, and the tremendous amount of labor that would be required to shange them (not to mention the coordination problems), the protogol has been highly resistant to further modification. However, there has been one set of ad hoc additions to the protosol made by many of the Hosta, notably all of the TIPS and TENEXes, single the additions were defined in early 1974.

The additions to the protogol were designed to deel with the lack of robustness in the flow dontrol mechanism, and the lack of a protogol mechanism for resynchronizing individual connections. The problem perceived to exist with the ingremental flow control mechanism is that it presumes a completely error-free transmission medium. However, low frequency Host software bugs, intermittent handwere bugs, and network fellures can cause an ALL control command to be lost, with the result that the buffer space accounting at sender and receiver gets out of synchronization. If this greates a situation in which the receiver believes buffers are eveilable, but the sender believes all buffers have been used, no further data can flow.

Use of the Host-Host RST (recet) command is inappropriate here, as it destroys all connections between the two Hosts. What is needed is a way to resynchronize only the affected connection without disturbing any others.

A second troublesome symptom of insonsistency in status information is the "half-closed" connection; after a service interruption or network partition, one NGP may believe that a donnection is still open, while the other believes that the connection is closed (does not exist). When such an

indonsistancy is discovered, the "open" and of the connection should be closed.

To echieve resynchronisation of apage ellocation, three new ad hoc control commands are defined; RAR (Reset Allocation by Receiver), RAS (Reset Allocation by Sender), and RAP (Reset Allocation Piesse), each of which takes one parameter; the link number for the connection.

The RAS command is sent from the Host sending on "link" to the Host receiving on "link". This command may be sent whenever the sending Host desires to resynchronize the status information essociated with the connection and doesn't have a message in transit through the network.

Some circumstances in which the sending Host may ghoose to do this are:

- After a timeout when there is traffic to move but no allocation (essumes that an allocation has been lost);
- 2) When an inconstitent event occurs associated with thet connection (e.g., an ALL is redelved which increases the message or bit space allocation beyond the allowed maximum).

- 3) After the sending host has suffered an interruption of network services
- 4) In response to a RAP (see below).

The RAR command is sent from the Host receiving on "link" to the Host sending on "link" in response to a RAS. It marks the completion of the connection resymphronization. When the RAR is returned the connection is in the known state of having no messages in transit and the allocations set to zero. The receiving Host may than start with a new allocation and normal message transmission can proceed. Since the RAR may be sent only in response to a RAS, there are no races in the resymphronization. All of the initiative lies with the sending Host.

If the receiving Hast detects an anomalous situation, however, a way is needed to inform the sending Host that a resynchronization is desirable. For this purpose, the RAP gommand is provided. It constitutes a "sudgestion" on the part of the receiving Host that the sending Host resynchronization sending Host is free to hence it or not as it sees fit. Since there is no obligatory response to a RAP, the regeiving Host may send them as frequently as it chooses and no harm can occur. For

example, if a message in excess of the allocate arrives, the receiving Host might send RAPs every few seconds until the sending Host replies, with no fears of races if one or more RAPs pass a RAS in the network,

An interruption of communication may result from a pertitioning of the subnet or from a service interruption on one of the communicating Hosts. It is undestrable to the up resources indefinitely under such direumstances, so the user is provided with the option of freeing up these resources (insluding himself) by unitaterally dispolving the connection. Here funitaterally means sending the CLS commend and slosing the connection without receiving the CLS acknowledgement.

When service is restored after such an interruption, the status information at the two ands of the connection may be out of synchronisation. One and believes that the connection is open and may proceed to use the connection. The disconnection end believes that the connection is closed (does not exist), and may proceed to reministrative communication by opening a new connection using the same socket pair or same link.

The resaymehronization needed here is to properly close the open and of the connection when the inconsistency is detacted.

This can be accomplished by specifying consistency checks and adding a new pair of commands. The two new ad hog commands are NXR (Non-existent Receive link) and NXS (Non-existent Send link), each of which takes as a parameter the link number.

The "missing CLS" situation described above sen manifest itself in two ways. The first way involves action taken by the NCP at the "epen" and of the connection. It may continue to send regular messages on the link of the helf-closed connection, or control messages referencing its link. The closed and should respond with an NXS if the message referred to a non-existent transmit link (e.g., was an ALL) or NXR if the message referred to a non-existent to a non-existent receive link (e.g., a data message). On receipt of such an NXS or NXR message, the NCP at the "open" and should close the connection by medifying its tables (without sending any CLS commend) thereby bringing both ends (nto agreement.

A second way this inconsistency can show up involves actions initiated by the NCP at the "sloged" end. It may (thinking the connection is closed) send an RFC to open a connection using the same socket or link. The NCP at the "open" end should detect the inconsistency when it receives such an RTS or STR command. In

this case, the NCP at the "open" end should close the connection (without sending any CLS command) to bring the two ends (nto agreement before responding to the RFC.

2.5.1.4 Constusions

Since the definition of the querent HosteHost protocol in early 1971, literally billions of bits have been successfully transferred between dissimilar ARPANET Hosts under its rules. Nevertheless, a number of problem areas have been identified during the intervening years, and subsequent protocol work in other ARPA programs, in other research networks, and in standards organizations have attempted to find solutions to many of these problems. There have also occasionally beam suggestions to nevise the ARPANET HosteHost protogol, but because of the large number of implementations which would be affected these have never met with widespread anthusiasm.

Some of the problem areas have been disquased in the previous paragraphs; these include the decision to make sonnections simplex (because of a "scardity" of link numbers); the lack of pobustness in the instrumental allocation agheme, which has been addressed by some Hosts with an ad hoc synchronization scheme, and the need to resynchronize NCP data

for single connections. The discussion of the Initial Connection Protocol in Section 2.5.2 will illustrate the class of difficulties which resulted from the decision to allow only one connection to a socket, in order to facilitate dynamic reconnection. Some of the other shortcomings are mentioned below.

The essumption that the ARPANET subnetwork provides a completely enrorefree transmission medium (i.e., not that it would always accurately return a RFNM or an Incomplete Transmission response) led to the design of a protosol that is not robust. Even if the assumption were true, there are still accessional errors in the IMP/Host interface and in the Host software which require corrections. There is not any requirement for extensive error detection and correction at the Host level, but errors will occur and the protocol should be able to survive them, As discussed in the previous section, it is frequently unable to do so.

On the other hand, some of the more esotaric expected communication scenarios have never actually occurred, and the protocol complexities designed to support them thus have some unused. In particular, the network usage has involved little or

no dynamic reconnection, and thus the protocol restrictions designed to make it possible have proven unnecessary. It is probable that one reason why the intense initial interest in resonnection has not given birth to its embediment in working systems is at least partly an adequating/authentication problem; no Host wants to provide service unless its accounting system can tell where to send the bill. This generally means, for existing Hosts, engaging in a dialog with the "logger" process. Development of a network-wide accounting system appears likely to be a prorequisite for dynamic movement of the logue of processing from Host to Hest.

Another unimplemented complexity is the type of dynamic buffer allocation anticipation by the GVS/RET mechanism. In general, large Hosts have been able to dedicate enough buffering to each connection so that GVSs are unnecessary. Conversely, the smallest Hosts, which might best profit from dynamic buffer allocation, have found that the memory space (and processing time) for the code needed to operate a dynamic buffering acheme is greater than the anticipated average gain. In addition, foverbooking schemes must be able to generate predictions of ever-allocation soon enough to account for the transmission and processing delay which will be inquired by a GVS request. Thus,

each Host contains the gode to progess a GVB, but none has the gode to generate one.

A major difficulty with the Host-Host protocol (s (te relience on the IMP/Host RFNM as a surrogate end-to-end adknowledgement. First, this reliance violates what subsequent workers in the protocol field have called "the principle of thin-wire design". The point is that if a protocol is designed to work over a "thin wire", it can be adapted to other sommunications medie, but if the protocol relies on the generation of auxiliary information (e.g., the RFNM) within the communication medium then it cannot sessity be adapted to use a different medium. This locks the Host into use of a particular technology and implementation, an undesirable situation both technology and economically.

As previously mentioned, the subnetwork responses which the Host needed (RFNM or Incomplete Transmission) were uniquely identified by the link, Since the IMPs initially prohibited more than one message outstending with a given link number, this restriction was carried into the Host protocol, However, when the IMPs discontinued this prohibition, the Hosts use of the link number as the exclusive identifier of which connection a

message belongs to made it impossible to relax the HosteHost protocol restriction. To illustrate this assertion, suppose that Heat A sends two messages on link n to Host B. At some later time a RFNM(n) and an Incomplete(n) arrive at A, and a single message on link m arrives at B. Host A tennet decide which of the two messages should be retransmitted, and Host 8 sannot deside if the message Just received immediately follows previous messages for the semmestion or if there is a gap. The Host-Host protocol restriction to a single outstanding message connection), as previously noted, restricts typical Host to Host bandwidth to about 10 Kbs, while the subnetwork can frequently deliver up to about 40 Kbs.

Chapter III

Another difficulty, elready briefly mentioned in Section 2.5.1.2.2. is the feat that every NCP implementar is faced with the necessity for "timing out" gertain protogol states in order to free valuable resources when a correspondent Host seems to be violeting protogol, but the timeouts are not uniform between Hosts. This situation sometimes leads to confusion when 1 40 heavily loaded Hosts are attempting to sommunicate, with one *timing out* states, and the other not doing so.

A supprising omission in the design of a protocol within a research community is the lack of specification of any instrumentation of Network Control Programs. It is also somewhat supprising that almost no NCPs have been instrumented, even if no instrumentation was specified. But in fact it is generally impossible to obtain answers to questions such as:

- What are the statistics about the number of connections open at one time?
- What are the statistiss about the number of messages sent or redelved over a sonnestion?
- What amounts of bit and message space are typically sliceted, and how is it used up?
- What is the ratio of control messages to data messages?

Pinelly, a very significant characteristic of the Hoste-Host protocol is that it is based on the concept of logical line switching in contrast to logical message switching. This has profound implications for its simplicity, flexibility, and rebustness. A way to understand these implications is to consider the alternative message-switching approach to Host-Host protocol design. As an aid to visualizing the approach, one of

the several specific proposals which has been made as an ARPANET alternative protocol is skatched briefly here.

Suppose that two processes, such as P and Q in Figure 2.5-2, wish to communicate; in particular, suppose P has a message to send to Q. Rather than setting up a connection (and its attendent flow control mechanism) with which to pass messages from P to Q, for such message Q is willing to receive, it sends a RECEIVE control message to Host A. When P has a message to send, it sends the data to its NCP. In whichever order these messages (RECEIVE or data) reach the Host A NCP, they eventually rendervous, and, at this time, the data is sent to Host B and the relevant entries in the table in the Host A NCP are discred. When the data gets to Host B, it matches a table entry left when the RECEIVE passed out of Host B. The table entries in Host B are then discred and the data is passed to process Q.

The significant features of this is that connections are not maintained over a sequence of messages but instead are set up and expire on a message-by-message basis. There are several advantages to this approachs

1. There is no need for a large number of intereNCP control messages to set up and break connections and donsequently no need for a special control channel.

Section 2 Design and Implementation

Figure 2,5e2 An Alternative messagesswitching approach

- 2. The RECEIVING progess has the opportunity after each message to stop the flow of messages. Further, since the RECEIVE process itself can be required to provide a buffer, the NCP is relieved of the tesk of providing a large buffering appealty.
- 3. Because connections exist only fleetingly, relatively dempies operations such as dynamically switching formections between a variety of processes are easy.
- 4. Errors have a minimal effect. For (natanse, (f a RECEIVE message is lost, the RECEIVING NCP will time it out and the process den do enother RECEIVE just as is normally done when the sender declines to send in response to a RECEIVE.
- This messege-switching protogol is suitable for implementation on small computers as well as large computers owing to its simplicity and the fact that there is no need for a large buffering deposity in the monitor.

2.5.2 Initial Connection Protocol

The Initial Connection Protocol (ICP) is a procedural protocol, rather than a formatting protocol, and applies only at the instant when a user attempts to make contact with a general service such as the "monitor" (or operating system) of a time-sharing terminal-oriented system.

The ICP owes its existence to the protogol restriction that permits a morket to be involved in only one connection at a time, coupled with the necessity of being able to "edvertise" the socket number at which a service (such as entry to a time-sharing system) may be found. Once the first user establishes a connection to the advertised socket, subsequent users will have to be turned away until the first user breaks the connection. Of course, the service-providing Host gould advertise a long list of socket numbers, and allow a user to work down the list until a free socket was found, but if the service-providing Host hopes to derive some benefit (e.g., income or popularity) from each user, turning a potential user away may represent lost opportunity.

An analogy exists with the situation feeing a business organization which makes sales via incoming telephone cells. The organization typically advertises a single telephone number.

However, the organization propures many incoming telephone direction, together with some type of automatic "ring up" equipment, so that if a sell is made to a number which is already in use, the cell will be switched (in a way invisible to the delier) to a free line. Naturally, at some point the capacity of the organization to handle more delies will be exceeded, but only at this point will the celler receive a donnection refusal (i,e., busy signel).

The ICP provides a function similar, sixhough not identical, to the "ring up" device. The switching sennot be sompletely invisible to the salling Host, singe that Host must know the setuel socket number involved in the sonnestion. Therefore the ICP specifies the following procedure to be followed by the server (salled) and user (salling) Hosts (refer to Figure 2.5-3).

A server progess makes evaluable a well-education dendesocket L and listens for a user (step 51). A user process initiates connection to sendesocket L from its receivement to specifying link "e" (U1). When the server receives en RTS ot socket L, it confirms the connection with an STR from sendesocket L to receivementally with a byte size of 32 bits (S2). The server then welts for an allocation from the user (S3). When the

SERVER		MAER	
811	Listen on socket L	Uli	RTS (U, L, +)
821	STR (L, U, 32)	U21	West for metch
831	Wait for allocation	451	ALL (a, m, b)
841	Send as deta the eccket number S	U4.	Receive the socket number 8 as date
351	CL8 (L. U)	U5:	CLS (U, L)
86;	RTS (8, U+3, x)	U6;	STR (U+3, 8, 61)
871	STR (8+1, U+2, n1)	U7 \$	RTS (U+2, 8+1, y)

Figure 2.5=31 Steps in Initial Connection Protocol

user receives confirmation of the connection (U2), it ellocates m (messages) and b (bits) for the connection using link MaM (U3), The server receives the allocation and sends, as regular data over the connection, the socket number 8, which occupies one 32-bit connection byte (84). The server sen then initiate the elosing of the connection between L and U (85). Once the user system has received the date message conteining 3 (U4) it can confirm (or initiate) the closing of the connection between L and U (U5).

The essence of the portion of the ICF described so far is to establish contect between the user and the edvertised socket just long enough for the server to select an arbitrary free socket and convey its number to the user, then break the connection. Since this should be able to happen rapidly, a small amount of queueing in the server should be sufficient to hold any dells which come in for the edvertised socket from other users during the connection's lifetime. The user has also received some assurance that the server is interested in establishing contact, as evidenced by the willingness of the server to provide a "caliback" socket.

Now, concurrently with the closing of the connection between L and U, the user and server exchange commends opening the connection from socket U+3 at the user to socket 8 at the server using link *x* and byte size ni (86 and U6). Also concurrently, commands are exchanged opening a connection from server socket 8+1 to user socket U+2 (87 and U7).

As can be seen from the above disquestion, the Initial Connection Protodol is relatively somplex. On the other hand, implementations are typically not very large, and tens or hundreds of thousands of connections have been established using

its it plearly is not too complex to be useful. However, most modern Hoste-Host protocol designs obviete the need for an ICP by requiring only one end of a connection to be unique. In view of the experience in ARPANET with dynamic reconnection, namely that there has been little or none, this is probably the proper course for protocol development.

2.5.3 TELNET Protocol

The TELNET protocol deals with the method of derrying out terminal-to-Host (or user to application program) gammunication, TELNET is derived from the phrase "Telecommunications Network", Since most of the Hosts initially connected to ARPANET provided some type of interestive, terminal-entiented service, there was intense interest in defining a protocol which would allow a human user, sitting at a keyboard/display (CRT or gharacter printer), to access a remote Host as though the terminal and the Host were directly connected; TELNET is the result.

2.5.3.1 Basic Concepts

The original TELNET protogol is built around the ideas of a Network Virtual Terminal (NVT), a dichotomy between users and servers, and a set of TELNET control signals.

The condept of a Network Virtual Terminal is an attempt to reduce the mapping problems refeed by interconnecting a veriety of terminals to a veriety of Hosts. Each different type of Host typically has a different internal character set and operating procedures. Similarly, each different terminal (usually even figureatible) has its own character set and/or timing

and operating constraints. If there are *H* Hosts and "I" terminals, and they are all directly connected to each other, then there are H*I total connections, and hands H*I total mappings required.

course, in the ARPANET terminals are not connected directly to packet ewitchess terminals are sonnected to Hosts (Sometimes mini=Hosts such as the terminal=support portion of the TIP, but Hosts nonetheless). Thus, a terminal gains access to a remote Host via the mediating intelligence of a logal Host The Virtual Terminal concept has the local process map DFOCESE. from the terminal-specific characteristics into e standard. Each Host is then required to implement only the mapping from its own internal characteristics to NVT characteristics. Thus the total number of mappings to be implemented is H+T, rather than the H=T required for direct mapping.

A second basic concept of the TELNET protocol is a dichotomy between "users" end "servers". For exemple, in a communication between a human at terminal (working through a local terminal support program) to a remote service process, only the human would have a preference as to where echoing should be done, or

would want to generate an interrupt. This viewpoint led to design decisions which made it difficult to use TELNET to support terminal-to-terminal, or process-te-process, communication.

Finally, it appeared necessary to have the user's (local) terminal support process and the server's operating system exchange some control information pertaining to the operating of the TELNET gonnection. The decision made was to embed TELNET control in the data atream to which it pertained. A small set of TELNET control codes was defined for this purpose.

2.5.3.2 Original Specification

TELNET was initially defined in the apring of 1971, elthough mugh of the definition was conveyed from implementer to implementer as "folklore" until a complete specification was dommitted to paper elmost a year later. However, this informality did not prevent TELNET from being rapidly and dompatibly implemented at a wide variety of sites.

2.5.3.2.1 Network Virtual Terminal Organization

The Network Virtual Terminal (NVT) is a bindirectional character device. The NVT has no timing sharesteristics. The characters are represented by 5 bit godes. The character codes 0

through 127 are the USASCII godes (code values are given in decimal). The codes 128 through 255 are used for special control signals. The NVT is described as having a printer and a keyboard. The printer responds to incoming data and the keyboard produces outgoing data.

2.5.3.2.2 NVT Printer

The NVT printer has an unspecified carriage width. The printer can produce representations of \$11.95 USASCII graphics. Of the 33 USASCII control codes, 8 have specific magning to the NVT printer, as shown in Figure 2.5=4. The remaining USASCII codes do not cause the NVT printer to take any action.

2.5.3.2.3 NVT Data Kayboard

The NVT Keyboard has keys or key sombinations or key sequences for generating all of the 128 USASCII codes. Note that although there are codes which have no effect on the NVT printer, the NVT Keyboard is capable of generating these codes.

The use of a standard, network-wide, intermediate representation of terminal code between sites is intended to eliminate the need for using and serving sites to keep

NAME	NE♥#TØ₫
NULL (NUL)	A no operation,
BELL (BEL)	Produces en audible or visible signal.
Back Space (88)	Backspaces the printer one character position.
Horizontal Tab (HT)	Moves the printer to next horizontal tab
Line Feed (LF)	Moves the printer to maxt line (keeping the same horizontal position).
Vertical Teb (VT)	Moves the printer to the next verticel tab stops
Form Feed (FF)	Moves the printer to the top of the next page:
Carriage Return (CR)	Moves the printer to the jeft margin of the current line.

Figure 2,5=4: NVT Printer Controls

information about the characteristics of each others terminals and terminal handling conventions. This approach can be successful, but only if the user, the using site, and the serving site assume certain responsibilities.

I. The serving site must specify how the intermediate gode will be mapped by it into the terminal godes that are expected at that site.

- 2. The user must be femiliar with that mapping.
- 3. The using site must provide some means for the user to enter all of the intermediate dodes, and as described below, special TELNET signals, as well as specify for the user how the signals from the serving site will be presented at the user terminal.

Since it is not known how the sites will specify the mepping between the network-wide standard code (7 bit ASCII in an 8 bit field) and the ecdes expected from their own terminals, it is necessary to permit the user to cause transmission of every one of the 128 ASCII codes.

2.5.3.2.4 Endeofeline Convention

The representation of the end of a physical line at a terminal is implemented differently on different network. Hosts, for example, some use a return (or new line) keys the terminal hardware both returns the garriage or printer to start of line and feeds the paper to the next line. In other implementations, the user hits carriage return and the hardware returns the derriage while the software sends the terminal a line feed. The network-wide representation is carriage return followed by line

feed. It represents the physical formatting that is being attempted, and is to be interpreted and appropriately translated by both using site and serving site.

Although TELNET defines the end of a line to be indicated by the ASCII character peir CR LF, several of the reel davides in the world have enty a single new line (NL) function. Several of computer systems have in some progress used the CR and LF functions to have semantic meaning larger than the format effect they provide, Further, several computer systems allow the CR and LF functions to be used separately (e.g., such that a line may be overprinted). One problem, for those TELNET (user) programs required to map the NVT into a device which only has a NL function, is how is the CR LF to be dealt with. A solution is to examine the character following the CR. If a LF is found, then perform the NL function; if enything else is found then back space to the beginning of the line. Another problem is the dase of a computer system which locally uses period, "." to cause new line function and which uses, in some programs, CR and LF for semantically significant operations. Suppose the user TEUNET sends the sequence CR LF. Does this mean "new line" on the "CR operation" followed by the "LF operation"? The solution to this problem is to require that TELNET programs send a CR which is not intended to be part of a CR LF pair as a CR NUL pair. Then the receiving program can always hold a CR and examine the next character to determine if a new line function is intended.

2.5.3.2.5 Break 8 | gmal

There is a special sentrol signed on some terminals that has no corresponding bit pattern in ASCII, but is transmitted by a special electric signal. This control signal is Attn on a 2741 and Break on a Teletype, Some systems treat the Break as an extre code available for use in conjunction with the data stream. For example, one system uses Break as a special aditing code meaning "delete the current line to this point," For this reason, the user must be provided with a way of generating this 129-th code from the NVT Keyboard, A representation of the Break code must also be specified for inclusion in the date streem; the TELNET control code "BREAK" (described in Section 2,5,3,2,8) is the defined signal.

2,5,3,2,6 Interrupt Signal

The user at a terminal connected to an interactive system generally has the capability of instructing the system to begin execution of user-specified programs. Since such programs are

result in unending execution loads, the user is normally provided with the capability of generating an "attention signal" which instructs the system to terminate the execution of the current process. (Some systems use the Break or Attn Key in this way, rather than as the 129-th gode described in the preceding perspace.)

One of the functions performed by a terminal gontrol program within an operating system is the scanning of an input streem for attention characters intended to stop an errent process and to return control to the executive. Terminal control programs which buffer input sometimes run out of space. When this happens to a local terminal's input streem, a "ball" or a question mark is echoed and the overflow character discarded, efter checking to see if it is the attention character. This strategy works wall in prectice, but it depends rather strongly on the intelligence of the human user, the inverient time delay in the input transmission system, and a lack of buffering between type-in and attention checking.

None of these conditions exists for interective traffic over the ARPANET. The empting Host connot control the speed (except

to slow it down) or the buffering within the using Most, nor can it even know whether a human user is supplying the input. It is thus recessory that the terminal control program or server TELNET not, in general, discard characters from a network input stream; instead it must suspend its acceptance of characters via the Host-Host flow central mechanism. Singe a Host may only send messages when there 1. room at the destinetion. dealing responsibility. for with too much input 1. thus transferred back to the using Host. This assures that no characters accepted by the using Host are inadvertently lost.

However, if the process in the serving Host stops accepting input, the pipeline of buffers between the user TELNET and remote process will fill up so that attention characters cannot get through to the serving executive. The solution to this problem dells for the user TELNET to send, on user-specified request, a Hostwhost interrupt signal foreing the server TELNET to switch input modes to process network input for estention characters in its network input, even if some input must be discarded while doing so. The effect of the interrupt signal to a server TELNET from its user is to cause the buffers between them to be emptied for the priority processing of extention characters. Thus the

user must be provided with a way of generating the "send=en-interrupt" signal from the NVT Keyboard. It must be noted that attention characters are Host specific and may be any of the 129 characters (128 ASCII characters plus Break); the requirements that the user must be able to generate the Break signal and the interrupt signal are independent.

To flip an attention againing server TELNET back into its normal mode, a special TELNET synchronization character (DATAMARK) is defined (see Section 2.5.3.2.8). When the server TELNET encounters this character, it returns to the strategy of eccepting terminal input only as buffer apage permits. It would not do to use the Host-Host signal along in place of the signal-DATAMARK combination in attention processing, because the position of the DATAMARK character in the TELNET input stream is required to determine where ettention processing ands and where normal mode input processing beings.

2,5,3,2,7 Local Functions

The ability of the user to dauge the using site TELNET to send any combination of ASCII characters in a string, and only that combination, is viewed as important to the user utility of the TELNET ASCII conventions. Because of this, some user sites

may find it necessary to implement a few additional function keys as part of the Iddal NVT implementation. These function keys are intended to provide local control between the terminal user and the terminal handling process. Two function keys which are likely to be universally required are:

Transmit Now: Transmit all date entered and locally buffered now:

Intended to be used with line mode operation:

Suppress end=of=line; Transmit all data entered and locally buffered now, and do not transmit the end=of=line immediately following this signal.

2.5.3.2.8 Control Codes

For control of the interaction between the user and server TELNET processes, a set of TELNET control codes are defined; these control codes are to be embedded in the TELNET data stream at appropriate points. In order to evoid interference with the normal data, the control codes are chosen from the gode values 128=255. Eight control codes are defined as described below:

BREAK; This is the 129mth data sode which the user may gause to be transmitted, as described in Section 2,5,3,2,5.

- NOP: This is a byte which may be used as a filler in the date stream to sligh text into a buffer location more convenient for the sender.
- TRANSPARENT and ESCRIC: These godes specify the escape to an alternate data_and_control code set. Since the control code set is also changed, there is no defined way to return to the normal ASCII data and control gode set described here.
- DATAMARK: This code is used in conjunction with the Host-Host protocol INE command to flip a TELNET server into and out of a special attention character againing mode as noted in Section 2.5.3.2.6.
- NOECHO: When this code is sent from user to server, it sake the server not to send echos of the trensmitted data. This is the default (start-up) condition. When this code is sent from server to user, it states that the server is not sending echos of the transmitted data. The server is permitted to send this code only as a reply to ECHO or NGECHO from the user, or to terminate the effect of HYI (see below). The server is prohibited from spontaneously sending NGECHO (or ECHO we see below) to evoid verious noce conditions which can develop if both user and server both try to select an echo mode simultaneously.

ECHO? When this code is sent from user to server, it asks the server to send echos of the transmited data. When this code is sent from server to user, it states that the server is sending echos of the transmitted data. The server is permitted to send this code only as a reply to ECHO or NDECHO.

HYI (Hide Your Input): The intention is that a server will send this code to a user system which is echoing locally (to the user) when the user is about to type something secret (e.g., a password). In this case, the user system is to suppress local schoing or everprint the input until the server sends a NOECHO code. In situations where the user system is not echoing locally, this code must not be sent by the server.

The HYI code presents some difficulty in that it is unclose how much is to be hidden. The server site usually knows how long the secret is but the user TELNET in general does not. Furthermore, if the user site dannot suppress the local echoing, there is a difficult implementation problem. One possibility is for the using site to overprint a full line with a mask, then have the user type his secret on the mask. If the secret were longer than one line, the use of the mask should be repeated.

The use of HYI might be avoided by having the serving site send a mask on which the user is to type the agenet information; This also presents a difficulty, however, in that the user system may be performing a mapping from terminal gode to NVT code which maps multiple keystrokes (such with a local echo) into single NVT characters. In addition, the user system may provide editing functions that map a deleted sharegter into several (rather than zero) characters. Thus, if the server sends a mask which it believes to be just long enough, some typing will be unmasked. On the other hand, sending a mask which is longer than the secret as expressed in NVT code may hide data which the user wants to see. Of course, the issue of masking arises only for those terminals whose internal echoing cannot be turned off.

2,5,3,3 Subsequent Modifications

By early 1973 there was a considerable amount of discontent with the TELNET protocol. There were several underlying themes to the discontent.

First, the dichotomy between "user" and "server" which we reflected in the details of many TELNET control structures (e.g., echoing, hiding input, interrupting) made it impossible to take advantage of the universality of TELNET implementations in using

TELNET for terminalwiseterminal or progesswiseprocess communication.

Second, there was continuing pressure for the inclusion of new controls in the data stream, but each new control gade which might be edded would require handling by each implementation. This had the effect of requiring unanimous "approval" of any addition; proposed additions never were able to muster unanimous support. As one example, the ISM 2741 terminal contains wired logic which gives control of the print head to sither the keyboard or the computer. The Keyboard is physically locked when the computer has control, so that the user cannot intersperse his input with the computer's output. Sings the source of output, the serving computer, knows when it is through sending output it would be useful if it gould send a control code to the user computer at that point, so the user computer sould relinquish control of the print head to the keyboard.

Third, the handling of where to generate echos was made entirely the responsibility of the user. However, each server is optimized to best handle terminals which either do their own echoing or do not, but not both. Therefore, the echoing conventions, which prohibit the server from initiating a change

in eaho made, seemed overly confining. The servers are burdened with users who are in the "wrong" mode, in which they might not have to be, and users, both human and machine, are burdened with remembering the proper echoing made, and explicitly setting it up, for all the different servers. Another echoing fasue was the desire of users, sepecially those whose BREWORK communication involved satellite links (e.g., from Hawaii or Europe), to have a highly interactive echoing mode in which the server system instructed the user system to eaho and not eaho classes of characters, with the definition of which classes to echo changing feirly dynamidally. This type of echoing can avoid long delays for every character while retaining the best human angineering features of the server systems which work in a character-at-a-time mode.

Fourth, it was recognized that almost every server system, as part of its own terminal support facilities, provided simple functions such as delete the last character, delete the last line, ring a bell to show you're still working, and so on. Unfortunately, almost every system chose a different user input string to activate each function. Both to simplify the task of a user in switching from system to system, as well as to facilitate the activation of such functions in process-to-process

communication, it seemed eppropriate to define standard estivation signals as part of the NVT.

For reasons of the kind described above, a "new" TELNET was defined in mid=1973 to replace the previous TELNET. This "new" TELNET was declared to be the official version, and a timetable for implementation was established calling for use of "old" TELNET to be discontinued in early 1974. However, the number of sites involved, the reduction in ARPA funding for protocol implementation work on a widespread basis, and a general rejustance to tamper with the working system gaused this schedule to be extended indefinitely. Nevertheless, many of the ARPANET Hosts, including the TIPs, have implemented the new TELNET and operate it in parallel with the old version.

The new TELNET differs from the old in three significant ways, which are discussed in the following sections. These are the principle of option negotiation, an expanded view of the NVT, and a revised view of the embedded sontrol codes:

2,5,3,3,1 Option Negotiation

The principle of negotiated options takes cognizance of the fact that many sites will wish to provide additional services

over and above those available within an NVT, and many users will have sophisticated terminals and would like to have elegant, rather than minimal, service. Independent of, but structured within, the TELNET Protocol verious "options" are sanctioned which can be used with the "DO, DON"T, WILL, WON"T" structure (discussed below) to allow a user and server to agree to use a more elaborate (or perhaps just different) set of conventions for their TELNET connection. Such options could include changing the character set, the scho mode, the line width, the page length, etc.

The basic strategy for setting up the use of options is to have either party (or both) initiate a request that some option take effect. The other party may then either accept or reject the request. If the request is accepted the option immediately takes effect; if it is rejected the essociated aspect of the connection remains as specified for an NVT. Clearly, a party may always refuse a request to anable, and must never refuse a request to disable, some option since all parties must be prepared to support the NVT.

The syntax of option negotiation has been set up so that if both parties request an option simultaneously, each will see the other's request as the positive acknowledgment of its own.

The symmetry of the negotiation syntex can potentially lead to nonterminating acknowledgment loops with each party seeing the incoming commands not as acknowledgments but as new requests which must be acknowledged. To prevent such loops, the following rules prevails

- Parties may only request a thange in option status: (.e., a party may not send out a "request" marely to announce what mode it is in.
- If a party receives what appears to be a request to enter bì some mode it is already in, the request should NOT be ecknowledged,
- Whenever one party sends an option command to a second party, a) whether as a request or an agknowledgment, and use of the option will have any effect on the processing of the data being sent from the first party to the second, then the command must be inserted in the date streem at the point where it is desired that it take effect. Some time will elapse between the transmission of a request and the receipt of an acknowledgment, which may be negative. Thus, a site may wish to buffer data, after requesting an option, until learns whether the request is accepted or rejected, in order to hide the "uncertainty period" from the user,

It is possible for requests initiated by processes to simulate a nonterminating request loop if the process responds to a rejection by merely repredenting the option. To prevent such loops from occurring, rejected requests should not be repeated until Something changes. Operationally, this can mann the process is running a different program, or the user has given enother command, or whetever makes sense in the context of the given process and the given option. A good rule of thumb is that a repreduest should only occur as a result of subsequent information from the other and of the connection or when demanded by local human intervention.

Option designers need not feel sanstrained by the somewhat limited syntex evaluable for option negotiation. The intent of the simple syntex is to make it easy to have options by making it some pertiaular option requires a nigher negotiation structure than possible within "DO, DON*T, WILL, WON*T", the proper tank is to use "DO, DON*T, WILL, WON*T" to establish that both parties understand the option, and once this is accomplished a more exotic syntex can be used freely. For example, a party might send a request to siter (establish) line length; If it is accompted, then a different syntex can be used for actually

negotieting the line length; such a "subenegotiation" perhaps including fields for minimum allowable, maximum allowable and desired line lengths. The important concept is that such expanded negotiations should never begin until some prior (standard) negotiation has assablished that both perties are speable of persing the expanded syntax.

In summary, WILL XXX is sent, by either party, to indicate that perty's desire (offer) to begin performing option XXX, DO XXX and DONFT XXX being its positive and negative similarly, DO XXX is sent to indicate a desire acknowledgments; (request) that the other party (i.e., the recipient of the DO) begin performing option XXX, WILL XXX and WONFT XXX being the positive and negative acknowledgments. Since the NYT is what is left when no options are enabled, the DONFT and MONFT responses are quarenteed to leave the connection in a state which both ends ean handle. Thus, all Hosts may implement their TELNET processes to be totally unaware of options that are not supported, simply returning a rejection to (f.e., refusing) any option request that sennot be understood.

2,5.3.3.2 The Expanded NVT

As with the old TELNET specification, the NVT has a character printer and Keyboard. The date is 7-bit ASCII code in an Sebit field. Echos will not, in the default NVT cross the network. The user is to be provided with a way to generate all 125 ASCII codes, plus the Break code, and the Interrupt signal. The end of line convention is as specified for old TELNET. There are no timing considerations for an NVT.

There are also five new keys defined (below) for the new NVT keyboard which map into five new NVT codes. The spirit of these extra codes is that they should represent a nature! extension of the mapping that already must be done from "NVT" into "local". as the NVT data byte 184 should be mapped into whatever the logal code for "uppersase D" (s, so the EC sharacter should mapped into whatever the lose! "Erese Character" function is, Further, just as the mapp(ng for 174 is somewhat arbitrary in an that has no "vertical ber" character, the environment character may have a somewhat arbitrary mapping (or none at all) there is no local "Erace Line" facility, (Similarly for format effectors: if the terminal actually dose have a "Vertical table then the mapping for YT is abvious, and only when the termine) does not have a vertical tab should the effect of VT be unpredictable.) The new codes eres

- IP (Interrupt Process): This is the "ettention signal" described in Section 2.5.1.2.6 as motivation for sending the Hoste-Host protocol interrupt. The interpretation of this character is generally to suspend, interrupt, abort, or terminate the process to which the NVT is connected.
- AC (Abort Cutput): Allow the current process to (eppear to) run to completion, but do not send its dutput to the user. Also send en interpupt signal (a Host-Host protocol INS and a DATAMARK) to ellow the sleeping of buffere external to the system receiving the AC and the logistion of subsequent data which should not be cleaped.
- AYT (Are You There); Send beek to the NVT same visible (i.e., printable) evidence that the AYT was received.
- EC (Erase Character); The recipient should delete the last preseding undeleted character or "print position" from the data streem.
- EL (Erece Line): The recipient should delete characters from the deta streem back to, but not including, the last "CR LF" sequence sent over the TELNET connection.

The definition of the NVT is also expended by the addition of guidelines for the transmission of data. Although a TELNET connection through the network is intrinsically full duplex, the NVT is to be viewed as a helf-duplex device operating in a line-buffered mode. That is, unless and until options are negotiated to the contrary, the following default conditions pertain to the transmission of data over the TELNET connections

1) Insofer as the evel(ability of local buffer space permits, data should be accumulated in the Host where it is generated until a complete line of data is ready for transmission, or until some locally-defined explicit signal to transmit occurs. This signal dould be generated either by a process or by a human user.

The motivation for this rule is the high cost, to some Hosts, of processing network input interrupts, coupled with the default NVT specification that "echos" do not treverse the network. Thus, it is reasonable to buffer some amount of data at its source. Many systems take some processing action at the end of each input line (even line printers or card punches frequently tend to work this way), so the trensmission should be triggered at the end of a line; On the other hand, a user or process may sometimes find it

negerary or desirable to provide data which does not terminate at the end of a line; therefore there must be methods of locally signalling that all buffered data should be transmitted immediately.

When e process has completed sending data to an NVT printer and has no queued input from the NVT keyboard for further processing (i.e., when a process at one and of a TELNET connection cannot proceed without input from the other end), the process must transmit the TELNET Go Ahead (GA) command. This rule is not intended to require that the TELNET GA command be sent from a terminal at the end of each line, since server Hosts do not normally require a special signal (in addition to endeofeline or other locally-defined characters) in order to command processing. Rather, the TELNET GA is designed to help a user's local Host operate a physically half duplex terminal which has a "lockeble" keyboard such as the IBM 2741.

2,5,3,3,3 Revised Control Code Structure

All TELNET commands constat of at least a two byte sequences the "Interpret as Command" (IAC) escape character (sode 255) followed by the sode for the command. The commands dealing with

ARPANET Completion Report Chapter III

aption negotiation are three byte sequences, the third byte being the code for the option referenced. This format was chosen so that as more comprehensive use of the "data space" (s. mede. (by negotiations from the basic NVT) sollisions of data bytes with reserved dommand values will be minimized, all such collisions requiring the inconvenience, and inefficiency, of "escaping" the data bytes into the stream. With the current setsup, only the IAC need be doubled to be sent as detay and the other 255 sodes may be passed transparently.

The defined TELNET commands are shown in Figures 2.5.5 and 2.5=6. These codes and code sequences have the indicated meaning only when immediately preceded by an IAC.

2.5.3.3.4 Defined Options

Because a participant in a TELNET conversation can refuse an option without understanding what it is, there is no pressure to limit the number of defined options, nor is there any pressure, other then the desire to provide more features, to implement defined options. By mid=1977 the 21 options had been defined as shown in Figure 2.5-7.

NAME	METHING
IAC	Data byte 253
NOP	No operation
Data Hark	This should always be accompanied by an INS on the control link
Break	NVT 129-th character
Interrupt Process	The fungtion IP
Abort Output	The fungtion AQ
Are You There	The fundtion AYT
Erase Character	The function EC
Erase Line	The function EL
Go Ahead	The GA signe!

Figure 2,5m5: TELNET Commends (Except Option) Negotiation

2,5,3,4 Constusions

Old TELNET has proven eminently suitable for the gennection of full-duplex terminals to server systems. It has been used, but with some difficulty, for communication between processes or between terminals. It does not provide a sufficient control structure to handle a basic helf-duplex terminal, or to incorporate all of the optional facilities that subsets of users desire. New TELNET was designed to deal with these shortcomings,

N▼#E	MEANING
WILL (option code)	Indicates the sender desires to begin performing, or sender the sender is now performing, the indicated aption
WON'T (option gode)	Indigates the refuse to perform, or continue performing, the indigated option
DO (aption code)	Indicates the request that the other party perform, or confirmation that the sender is expecting the other party to perform, the indicated option
(ebes noitas) Tènda	Indicates the demand that the other party stop performing, or confirmation that the sender is no longer expecting the other party to perform, the indicated option
SB (option code)	Indicates that what follows is a list of subnegotiation parameters pertaining to the indicated option. The list will be terminated by IAC SE
SE	Submagotiation End = indicates the end of a list of option submagotiation parameters

Figure 2.5-6: TELNET Option Negotiation Commands

and appears to do so successfully. Nevertheless, the experience with TELNET illustrates the transmission inertia associated with a

- 1. Switch to binery transmission mode
- 2. Send Eahos over the TELNET connections
- 3. Reconnect the TELNET connection
- 4. Suppress transmission of GA (GowAhead)
- 5. Select a convenient message size
- 6. Provide surrentestatus=of=options information
- 7, Send a timingement (to indicate a roundetrip over the data connections)
- 8. Enter a complex echoing made based on classes of characters to be achoed at each end.
- 9, 3elect a convenient print line length
- 18. Select a convenient page length
- 11. Choose the best system to provide any necessary padding
- 12. Select horizontal tabatons
- 13. Choose the best system to properly handle horizontal tabs
- 14. Choose the best system to properly hendle formfeeds
- 15. Select vertical tabatons
- 16. Choose the best system to properly handle vertical tabs
- 17, Choose the best system to properly handle linefeeds
- 18. Agree to use a particular "extended ASCII" character set
- 19, Log the user off the server system
- 28. Use a particular textecompaction scheme
- 21. Switch from operation of an NVT to a "virtuel" Data Entry Terminal

Flaure 2.5-71 TELNET Options

program that works. The discarding of old TELNET is now about four years overdue (out of a total age of less than seven years). An important conclusion, therefore, is that the first version of a protocol which must be implemented by many different groups will persist in spite of later improvements; it is therefore very important to approach perfection as closely as possible from the beginning. This inertia can be contrasted with the relative ease with which the IMP software, built by a single group on only one

type of mechine, has been drastically changed during the course of the ARPANET project.

A second gonelysion to be drawn from the TELNET experience is the Unifying power of the "virtual device" concept, coupled with a mandatory minimal implementation extended by negotiable options. The virtual device concept pervades modern terminal protocol development throughout the world.

2.5.4 File Transfer Protogol (FTP)

The File Transfer Protoco! specifies a mechanism for the transfer of complete files (rether than, for example, groups of "records" retrieved from a file) between ARPANET Hosts, The primary function of FTP is to transfer files efficiently and reliably among Hosts and to allow the convenient use of remote file storage capabilities.

The objectives of FTP ares

- 1) to promote ehering of files (computer programs end/or data)
- 2) to enacurage indirect or implicit (via programs) use of remote computers
- 3) to shimld a user from variations in file storage systems among Hosts
- 4) to transfer data reliability and efficiently,

FTP, though usable directly by a user at a terminal, is designed mainly for use by programs. The attempt is to satisfy the diverse needs of users of maxisHosts, minisHosts, TIPs, and the Datacomputer, with a simple, and easily implemented protocol design.

The effort to develop a File Trensfer Protogol began with the condeptualization of two levels; a "deta trensfer protogol" to govern the transmission of new data, and a "file transfer protogol" which would control the data flow, establish the names and locations of files to be moved, user addess rights, and so on, preliminary specifications for these two levels were developed, but it eventually became alean that the two layers were ectually only the formatespecification and procedure-specification aspects of a single protogol, the current FTP.

2.5.4.1 Benio Concepts

As with TELNET, the FTP can be best understood in terms of a user and a server. Unlike TELNET, however, the protogol expects the user to be a process, rather than a human at a keyboard. Nevertheless, some care has been taken to make it possible, if not convenient, for the user to be a human.

Figure 2.5m8 provides a sketch of the overall model of the FTP system. The "Server FTP Process" gan be further subdivided into a server Protocol Interpreter (PI) which interfaces to the dominant/response path via TELNET, and a Date Transfer Process (DTP) which moves the data between the date connection and the

Figure 2,5=8: ARPANET File Transfer Protocol

file system. Similarly, the "user FTP Process" can be subdivided into a Protocol Interpreter, a Data Transfer module, and a User Interfece which mediates between the Protocol Interpreter and the local user support facilities.

In a typical section, the user would be expected to activate the User FTP Process. The user-protocol interpreter then initiates the TELNET connections. At the initiation of the user, standard FTP commands are generated by the user-PI and transmitted to the server process via the TELNET connections. (The user may establish a direct TELNET connection to the server-FTP, from a TIP terminal for example, and generate

standard FTP commands himself, byepassing the usereFTP process;)

Standard replies are gent from the serverePI to the userePI over the TELNET connections in response to the commands.

The FTP dommends specify the peremeters for the data sommection (representation type, file structure, end transfer mode), and the neture of file system operation (stone, retrievo, append, delete, etg.). The user=DTP or its designate should "listen" on the specified data socket, and the server initiate the data gonnection and data transfer in accordance with the specified parameters. It should be noted that the data socket need not be in the same Host that initiates the FTP dommands via the TELNET gonnections, but the user or his user=FTP process must ensure a "listen" on the specified data cooket. It should also be noted that two data connections, one for send and the other for receive, may exist simultaneously.

In enother mituation a year might wish to transfer files between two Hosts, neither of which is his local Host. He sets up TELNET connections to the two servers and then arranges for a date connection between them. In this manner control information is passed to the user-PI but date is transferred between the server date transfer processes.

FTP commands are "TELNET strings" term(nated by the "TELNET end of line code". The command codes themselves are alphabetic characters terminated by the character <SP> (Space) if parameters follow and CR=LP as an end-of-line indication if not.

Replies to File Trensfer Protogo! commands were devised to ensure the synchronization of requests and actions in the process of file transfer, and to guarantee that the user process always knows the state of the Server, Every command must generate at least one reply, although there may be more than one; in the latter case, the multiple replies must be easily distinguished. In addition, some commands occur in sequential groups. The replies show the existence of an intermediate state if all preceding gommands have been successful. A failure at any point in the sequence necessitates the repetition of the entire sequence from the beginning.

An FTP reply gensists of a three digit number (transmitted as three elphanumer(s characters) followed by some text. The number is intended for use by automata to determine what state to enter next; the text is intended for the human user. It is intended that the three digits contain enough encoded information that the User-PI will not need to examine the text and may either

discard it or pass it on to the user, as appropriate. In particular, the text may be server-dependent, so there are likely to be varying texts for each reply sede.

There will be eases where the text is longer than a single line. In these eases the complete text must be bracketed so the User-process knows when it may stop reading the reply (i.e. stop processing input on the TELNET connection) and go do other things. This requires a special format on the first line to indigate that more than one line is coming, and another on the last line to designate it as the last. At least one of these must contain the appropriate reply sade to indigate the state of the transaction. The format chosen for multipline replies which meets these constraints is that the first line begins with the exact required reply code, followed immediately by a hyphen, """, followed by text. The last line begins with the same code, followed immediately by space <\$P>, eptionally some text, and CR

123-First line
Second line
234 A line beginning with numbers
123 The lest line

The user-process then simply needs to search for the second occurrence of the same reply code, followed by 45Pb (Space), at the beginning of a line, and ignore all intermediary lines. If an intermediary line begins with a 3-digit number, the Server must pad the front to avoid confusion.

The three digits of the reply each have a special significance. This is intended to allow a range of very simple to very sophisticated response by the user-spresss. The first digit denotes whether the response is good, bad or insomplete. An unsophisticated user-spresses will be able to determine its next action (proceed as planned, redo, retrench, etc.) by simply examining this first digit. A user-spresses that wants to know approximately what kind of error occurred (e.g. file system error, command syntax error) may examine the second digit, reserving the third digit for the finest gradation of information.

There are five values for the first digit of the reply codes

Lyz Positive Preliminary reply: The requested action is being initiated; expect another reply before proceeding with a new dommand. This type of reply can be used to indicate that the command was accepted and the user=process may now pay

Section 2 Design and Implementation

attention to the data connections, for implementations where simultaneous monitoring (a difficult.

- 2 v z Positive Completion reply: The requested action has been successfully completed. A new request may be initiated.
- Svz Positive Intermediate replys The dommand bas bean accepted, but the requested action is being held in abeyance, pending receipt of jurther information. The user should send another command specifying this information, This reply is used in command sequence groups.
- 441 Transient Negative Completion reply: The sommand was not accepted and the requested action did not take place, but the error condition is temporary and the action may be requested again. The user should return to the beginning of the command sequence, if any. It is difficult to assign a meaning to "transient", particularly when two distinct sites (Server and Usereprocesses) have to 20766 on. the interpretation. Each reply in the 4yx datagory might have a slightly different time value, but the intent is that the user-process is engagraged to try again. A rule of thumb in determining if a reply fits into the 4vz or the 5vz (Permanent Negative) dategory is that replies are 4vz if the

commands can be repeated without any change in command form or in properties of the User or Server (e.g. the command is spelled the same with the same arguments used; the user does not change his file addess or user name; the server does not out up a new implementation.)

Permanent Negative Completion reply: The dommand was not accepted and the requested action did not take place. The User-process is discouraged from repeating the exact request (in the same sequence): Even some "permanent" error conditions can be corrected; so the human user may want to direct his User-process to reinitiate the command sequence by direct action at some point in the future (e.g. after the spelling has been changed, or the user has altered his directory status.)

The following function groupings are encoded in the second digits

- x0z Syntex; These replies refer to syntex errors, syntectically correct commands that don't fit any functional category, unimplemented or superfluous commands.
- xiz Information; These are replies to requests for information, such as status or help,

- x2z Connections: Replies referring to the TELNET and data connections.
- x3z Authentication and addountings Replies for the logon process and accounting procedures.
- x5s file system: These replies indicate the status of the Server file system vis=e=vis the requested transfer or ther file system ention.

The third digit gives a yet finer gradation of meaning in each of the function categories.

2.5.4.2 Original Specification

There are five areas of FTP specification. Three of them (representation type, file atructure, and transmission mode) deal with the data transfer espects of the protogol, while the other two (data error recovery, command structure) deal with the protogol aspects of the protogol. Each of these five areas is described briefly below.

2.5.4.2.1 Data Representation Types

Data is transferred from a storage device in the sending Host to a storage device in the receiving Host. Often it is

Section 2 Design and Implementation

necessary to perform certain transformations on the data because storage representations in the two systems are different. For example, NVT=AGCII has different data storage representations in different systems. PDP=10fs generally store NVT=ASCII as five 7-bit ASCII characters, leftwlustified in a 36-bit words 360's NVT-ASCII as 8-bit EBCDIC codes: Multica stores NVT-ASCII as four 9-bit characters in a 36-bit word. It may be desirable to convert characters into the standard NVT-ASCII representation when transmitting text between dissimilar systems. The sending reselving sites would have to perform the necessary and transformations between the standard representation and their internal representations.

different problem in representation arises transmitting binary data (not character godes) between systems with different word lengths. It is not always glear how the sender should send data, and the reseiver store it. example, when transmitting 32-bit bytes from a 32-bit word-longth svatem to a 36-bit word-length system, it may be desirable (for reasons of efficiency and usefulness) to store the 32-bit bytes rightulustified in a 36mbit word in the latter system. In any ease, the user should have the option of specifying data representation and transformation functions, FTP provides for

very limited data type representations; the following types are defined:

ASCII Format: This is the default type and must be accepted by all FTP implementations. It is intended primarily for the transfer of text files, except when both Hosts would find the EBCDIC type more donvenient. The sender converts the data from his internal character representation to the standard 5-bit NVT-ASCII representation as defined in the TELNET specification. The receiver will convert the data from the standard form to his own internal form. In accordance with the NVT standard, the <CRLF> sequence should be used, where necessary, to denote the end of a line of text.

EBCDIC Format: This type is intended for efficient transfer between Hosts which use EBCDIC for their internal character representation. For transmission the data are represented as 8-bit EBCDIC characters. The character gode is the only difference between the functional specifications of EBCDIC and ASCII types. End=of=line (as opposed to end=of=record) will probably be rarely used with EBCDIC type for purposes of denoting structure, but where it is necessary the <NL> character should be used.

A Character file may be transferred to a Host for one of three purposes: for printing, for storage and later retrieval, or for processing. If a file is sent for printing, the receiving Hast must know how the vertical format control is represented. In the second case, it must be possible to store a file at a Host and then retrieve it later in exactly the same form. Finally, it aught to be possible to move a file from one Host to another and process the file at the second Host without undue trouble. A single ASCII or ESCDIC formet does not satisfy all these conditions and so these types have a second parameter specifying one of the following three formate:

- Non-print: This is the default format. The file need contain no vertical format information. Normally, this format will be used with files destined for prosessing or just storage.
- TELNET Format Controles The file contains ASCII/EBCDIC vertical format controls (f.e., <CR>, <LF>, <NL>, <VT>, <FF>) which the printer process will interpret appropriately, <CRLF>, in exactly this sequence, also denotes end-of-line,
- Cappings Control (ASA): The file contains ASA (FORTRAN) vertical format gontrol characters. In a line or a record, formatted according to the ASA Standard, the first character is not to

be printed. Instead it should be used to determine the vertical movement of the paper which should take place before the rest of the resord is printed.

Image: The data are sent as contiguous bits which the receiving after must store as contiguous bits. The structure of the storege system might necessitate the padding of the file (or of each record, for a record-structured file) to some convenient boundary (byte, word on block). This padding may occur only at the end of the file (or at the end of each record) and there must be a way of identifying the padding bits so that they may be stripped off if the file (a retrieved. The padding transformation should be well publicized to enable a user to process a file at the storage site. Image type is intended for the afficient storage and retrievel of files and for the transfer of binary data.

Of dourse, a file must be stored and retrieved with the same parameters if the retrieved version is to be (dentical to the version originally transmitted. Conversely, FIP implementations must return a file identical to the original if the parameters used to store and retrieve a file are the same.

2.5.4.2.2 File Structures

FTP recognizes two file structures; a *pupe* file structure with no internal substructure, and a record structure in which the file is composed of a series of records. Of source, these two recognized structures represent only a small percentage of the total number of file structures which have been defined, but they are sufficient to handle most of the files which ARPANET Hosts wish to exchange,

The "natural" structure of a file will depend on which Host stores the file. A source-code file will usually be stored on an IBM 368 in fixed length regards but on a PDP-18 as a stream of characters partitioned into lines, for example by «CRLF». If the transfer of files between such disperses sites is to be useful, there must be some way for one site to recognize the other?s essumptions about the file.

With some sites being naturally file-oriented and others naturally record-oriented there may be problems if a file with one structure is sent to a Host oriented to the other. If a taxt file is sent with record-structure to a Host which is file oriented, then that Host must apply an internal transformation to the file based on the record structure. Obviously this

transformation should be useful but it must also be invertible so that an identical file may be retrieved using record structure.

In the case of a file being sent with file-structure to a record-oriented Host, there exists the question of what priteris the Host should use to divide the file into records which can be processed lossily. If this division is necessary the FTP implementation should use the end-of-line sequence, <CRLF> for ASCII, or <NL> for ESCDIC, text files as the delimiter. If an FTP implementation adopts this texhnique, it must be prepared to reverse the transformation if the file is retrieved with file-structure.

2.5.4.2.3 Trenemission Modes

The finel consideration in transferring date is choosing the appropriate transmission mode. There are three modes defined in FTP1 one which formets the date and ellows for restart procedures; one which also compresses the date for efficient transfer; and one which passes the date with little or no processing. In this last case the mode interests with the structure attribute to determine the type of processing. In the compressed mode the representation type determines the filler byte.

All data transfers must be completed with an end-of-file (EOF) which may be explicitly stated or implied by the closing of the data connection. For files with record structure, all the end-of-record markers (EOR) are explicit, including the final one.

The following transmission modes are defined:

Stream: The data is transmitted as a stream of bytes. There is no restriction on the representation type used; record structures are allowed. In a record structured file EOR and EOF will each be indicated by a two-byte control code of whatever byte size is used for the transfer. The first byte of the centrol code will be all ones, the escape character. The second byte will have value 1 for EOR and value 2 for EOF. EOR end EOF may be indicated together on the last byte transmitted by the value 3. If a byte of all ones was intended to be sent as data, it should be repeated in the second byte of the control code.

If the file does not have record structure, the EOF is indicated by the sending Host signing the data connection and all bytes are data bytes.

For the purpose of standardized transfer, the sending Host translates its internal and of line or end of regard denotation into the representation prescribed by the transfer mode and file structure, and the receiving Host performs the inverse translation to its internal denotation. Since these transfermations imply extra work for some systems, identical systems transferring non-record structured text files might wish to use a binary representation and stream mode for the transfer.

Block! The file is transmitted as a series of data blocks preceded by one or more header bytes. The header bytes contain a count field, and descriptor code. The count field indicates the total length of the data block in bytes, thus marking the beginning of the next data block (there are no filter bits). The descriptor code defines: last block in the file (EOF), last block in the resond (EOR), restart marker (see Section 2,5,4,2,4), or suspect data (i.e., the data being transferred is suspected of errors and is not reliable). This last code is not intended for error control within FTP, It is motivated by the desire of sitos exchanging certain types of data (e.g., seismic or weather data) to send and receive all the data despite local errors

(such as "magnetic tape read errors"), but to indicate in the transmission that parts(n potions are suspect).

Record structures are allowed in this mode, and any representation type may be used. There is no restriction on the transfer byte size. With this anguading more than one descriptor coded condition may exist for a particular block. As many bits are necessary may be flagged.

The restert marker is embedded in the data stream as an integral number of 8-bit bytes representing printable characters in the language being used over the TELNET connection (e.g., default is NVT-ASCII).

Compressed: The file is transmitted as series of bytes of the size specified by the BYTE command, There are three kinds of information to be sent; regular data, sent in a byte string; compressed data, consisting of replications or filler; and control information, sent in a two-byte escape sequence. If n bytes of regular data are sent, these n bytes are preceded by a byte with the left-most bit set to 2 and the other bits containing the number n.

If n replications of the data byte d are to be transmitted, they can be compressed into two bytes. The left-most two bits of the first byte are set to 10 and the remaining bits contain the number n. The second byte contains the value d.

A string of m filler bytes can be compressed into a single byte, where the filler byte veries with the representation type. If the type is ASCII or ESCDIC the filler byte is space. If the type is Image or Local byte, the filler is a zero byte. The compressed filler is a single byte with the left-most two bits set to 11 and the remaining bits containing the number n.

The escape sequence is a double byte, the first of which is the escape byte (all zeroes) and the second of which contains descriptor codes as defined in Block mode. The descriptor sodes have the same meaning as in Block mode and apply to the succeeding string of bytes.

Compressed mode is useful for obtaining increased bandwidth on very large network transmissions at a little extra CPU cost. It is most afficient when the byte size chosen is that of the word size of the transmitting Host.

2.5.4.2.4 Error Recovery

There is no provision for detecting bits lost or scrembled in data transfer. However, a restart procedure is provided to protect users from gross system failures (including failures of a Host, an FTP=process, or the IMP subnet):

The restart procedure is defined only for the block and compressed modes of data trensfer. It requires the sender of data to insert a special marker code in the data stream with some marker information. The marker information has meaning only to the sender, but must consist of printable characters in the code of the TELNET connection. The marker could represent a bitmount, a recordedount, or any other information by which a system may identify a data checkpoint. The receiver of data, if implements the restart procedure, would then mark the corresponding position of this marker in the receiving system, and return this information to the user.

In the event of a system failure, the user can restart the data transfer by identifying the marker point with the FTP restart procedure. The following example illustrates the use of the restart procedure.

The sender of the data inserts an appropriate marker block in the data stream at a convenient point, The receiving Host marks the corresponding data point in its file system and conveys the last known sender and receiver marker information to the user, either directly or over the TELNET connection, depending on who is the sender. In the event of a system failure, the user or controller process resterts the server at the last server marker by sending a restart command with the server's marker code as its argument. The restart command is transmitted over the TELNET connection and is immediately followed by the command (such as Store or Retrieve) which was being executed when the system failure occurred.

2.5.4.2.5 FTP Commands

There are three general glasses of FTP commands. First, there are commands which deal with the issues of access control and accountings logging on and logging off the server system.

Second, there are the sommands which desi with the data transfer issues such as choosing a representation type, file structure, and transmission mode. Commands to establish socket numbers and a byte size for the data connection are also included. Defaults exist for all of the data transfer

parameters, so these commands are needed only if the default are to be changed.

Third, there are the service commends which define the file transfer or the file system function requested by the user. The ergument of an FTP service commend will normally be a pathname. The syntax of pathnames must conform to server site conventions (with standard defaults applicable), and the language conventions of the TELNET connection. The commends may be in any order except that a "rename from" commend must be followed by a "rename to" commend and the restart commend must be followed by the interrupted service commend. The data, when transferred in response to FTP service commends, is always sent over the data connection, except for certain informative replies. The commends to specify FTP service (nelydes

- Retrieve causes the server to send a gopy of the named file over the data connection.
- Store w causes the server to eccept the data transferred over the data connection and store it with the given pathname.
- Append w payers the server to addept the data transferred via the data connection and append it to the named file, or create the named file if it does not exist.

Allocate # reserve space for file storage,

- Restart = accept a restart marker (see Section 2.5,4.2.4) and prepare to resume a data transfer.
- Rename = a "from" and a "to" dommand, used in pairs, to rename a file stored at the server.
- Abort = eny action of the lastegiven unfinished command, especially the transfer of data.
- Delete = a file stored at the server,
- List a depending on parameters, send a list of all the files belonging to the user or stored in a named directory, or datails about the storage of a named file.
- Status = causes the server to send a status response over the TELNET connection; in general this is used to obtain the server's view of the status of a data transfer which is in progress.

2.5.4.3 Conclusions

As previously noted, the official FTP was developed fairly late relative to the HostaHost and TELNET protocols, and thus was able to incorporate features to deal with some of the shortcomings experienced during use of these protocols. The two most interesting ideas in the FTP are the effort to separate the

data and control functions so completely (i.e., with apparate connections using different formatting/encoding), and the structuring of replies so that they can be used either by human or automaton.

The separation of data and control was seen to be successful at the Host-Host levels it is for less successful in FTP, where synchronization problems (both over the network and within a server system) have forced some control to be embedded in the data stream, and have also in general prohibited the user system from having multiple commands outstanding.

The use of responses (and commands) which can be generated and interpreted both by humans and programs (even on radically different systems) is probably the major technical contribution of the ARPANET FTP to the body of protocol ideas. The concept has been proven to work quite well in practice, and has been copied in the development of protocols for other networks.

On the other hand, the FTP has been oriticised for attempting to obtain generality by the method of including a little of something for everyone; witness the six representation types (ASCII and EBCDIC each with three subtypes), two file structures, and three transmission modes = a total of 36

combinations. It has been suggested by several protocol workers that a more retional approach is to define a "Network Virtual File System"; a single, general, simple, albeit inefficient protocol that sould be used by all casual users of file transfer and similar functions. This could even be embedded in TELNET; All communication between pairs (or other natural groups) of serious users should be done using special-purpose protocols. While this may represent an extreme view, it would have been much easier to implement and would have encouraged users with large volumes of data, and a desire for efficiency, to use a special-purpose conversion programs, well-ematched to the specific application, rather than a General-purpose conversion at each end,

Section 2 Design and Implementation

2.5.5 Other HighwLevel Protocole

There are a number of additional highelevel protocols, some "official" and others ad hoo, which have been designed for the ARPANET. This section contains brief descriptions of three such protocols: the official Remote Job Entry protocol, the "competing" ad hoo Remote Job Service protocol, and the official Graphics protocol.

2.5.5.1 Remote Job Entry Protocol

Remote job entry is the machenism whereby a user at one location causes a batcheprocessing job to be run at some other location. This protocol specifies the Network standard procedures for such a user to communicate over the Network with a remote batcheprocessing server, causing that server to retrieve a job-input file, process the job, and deliver the job's output file(s) to a remote location. The protocol uses a TELNET connection for all dentrol communication between the user and the server RJE processes. The serveresite then uses the File Transfer Protocol to retrieve the job-input file and to deliver the output file(s). This mechanism is illustrated in Figure 2.5-9. It should be noted that the RJE user need not be located at the Host which sends the input or receives the output (i.e., the Host containing the FTP server).

Figure 2,509: ARPANET Official RJE Protocol

There are two types of users: direct users (persons) and user processes. The direct user communicates from an interactive terminal attached to a Hest. This user may cause the input and/or output to be retrieved/sent on a specific socket at the specified Host (such as for card readers or printers on a TIP), or the user may have such files transferred by file-name using File Transfer Protocol. The other type of user is an RJE User-process in one remote Host dommunicating with the RJE Server-process in another Host. This type of user ultimately receives its instructions from a human user, but through some unspecified indirect means. The commend and response streems of this protocol are designed to be readily used and interpreted by both the human user and the user process.

A perticular user site may ghoose to establish the TELNET control connection for each logical job or may leave the control connection open for extended periods. If the control connection is left open, then mulitple job-files may be directed to be retrieved or optionally (to servers that are able to determine the end of each logical job and form several jobs out of one input file) one continuous retrieval may be done (as from a TIP serd reader). This then forms a "hot" cord reader to a particular server with the TELNET connection serving as a "job

monitor. Since the output is always transferred a job at a time per connection to the supput secket, the putput from this "hot" reader would appear when ready as if to a "hot" printer. Another possibility for more complex hosts is to attach an RJE User-process to a gard reader and take instructions from a lead control gard, gausing an RJE control TELNET to be opened to the appropriate Host with appropriate logson and input retrievel commands. This card reader would appear to the human user as a Network "hot" card reader.

The RJE commends and responses follow the same general rules as defined for FTP. In fact, the reply codes are primary those defined for FTP, with only same minor additions. The commend structure is somewhat more complicated than exists for FTP, for three reasons. First, the RJE server site contains an FTP user which must typically "log in" to the Host which is the source or destination or the Jeb data. Thus, the RJE user must not only supply accounting and authentisation data to gain entrance to the RJE server; the user must also supply authentisation and accounting information to be used by the RJE server to gain access to the files. Second, the user must be able to apacify a Host system from which the RJE server is to obtain the input and a system to which the RJE server is to send the output; these

avatems need not be the Host at which the user is located. need they be the same as each other. Thirds the user must specify whather the RJE server is to actively try to send the output after it becomes evaliable or wait for the user to retrieve it. The user must also specify whether the RJE server is to discard or retain the output after it has been sent to the destination Host.

Although the official RJE protogol was specified by late 1972, it has probably never been implemented. This is further testimony to the inertia of a working existing system. sompeting evatem in this case is the Network Remote Job Service protocol specified as an "interim" measure a year and a half earlier; the users end servers with a real need to perform remote batch computing implemented this protocol immediately and never bothered to change,

2.5.5.2 Network Remote Job Service

NETRUS is the protocol for the remote jeb entry service on the IBM 360 Model 91 at the UCLA Campus Computing Natwork (CCN), NETRJS ellows the user at a remote Host to eccess CCNfs RJS (MRemote Job Services) subsystem, which provides remote Job entry service to real remote batch (card reader/line printer) terminels over direct communications lines as well as to the ARPANET.

To use NETRJS, a user at a remote Host needs a NETRJS user process to communicate with one of the NETRJS server processes at CCN. Each active NETRJS user process appears to RJS as a Separate (v(rtual) remote batch terminal (VRST). A VRST may have virtual dard readers, printers, and punches. Through a virtual dard reader a user can transmit a stream of dard images comprising one or more 08/360 Jobs, complete with Job Control Language, to CCN. These lobs are specied into CCNPs batch system (08/366 MVT) and run according to their priority. RJS will automatically return the print and/or punch output images which are created by these jobs to the virtual printer and/or card punch at the VRBT from which the lob same (or to a different destination specified in the JCL). The remote user can wait for the output, or can sign off and sign back on later to receive it. Figure 2.5=10 illustrates the processes and connections involved in the use of NETRJS.

The VRBT is assumed to be under the control of the user's console; this serves the function of an RJS remote operator densole. To initiate a NETRJS session, the remote user must establish a TELNET connection and sign into RJS. Once signed in, he can use his console to issue domands to RJS and to receive status, confirmation, and arror messages from RJS. Different

Figure 2,5=10: NETRJS Protocol

VRBTs are distinguished by 8=character terminal idfs. There may be more than one VRBT using RJS simultaneously from the same remote Host.

When a VRBT starts a session, it has a choice of two ICP sockets, depending upon whether it is an ASCII or an EBCDIC virtual terminal. An EBCDIC virtual terminal transmits and received its data as transparent streams of 8 bit bytes (since ECN is an EBCDIC installation). It is expected that a user at an ASCII installation, however, will want his VRBT declared ABCII; RJS will then translate the input stream from ASCII to EBCDIC and translate the printer stream back to ASCII. This will allow the user to employ his local text editor for preparing input to CCN and for examining cutput. The punch stream will always be transparent, for outputting "binary decks". The choice of code for the operator console connections is independent of declared terminal type; in partiquiar, they always use ASCII under TELNET protocol, even from an EBCDIC VRBT.

NETRUS protodol provides data compression, replacing repeated blanks or other characters by repeat counts, However, when the terminal id is assigned by CCN, a particular network terminal may be specified as using no data compression. In this

dase, NETRJS will simply truncate trailing blanks and send records in a simple "op godewlangth=data" form, called "truncated format".

A job stream for submission to RJS at CCN is a series of logical records, each of which is a card image. A card image may be at most 80 characters long, to match the requirements of 08/360 for job input. The user can submit a "stack" of successive jobs through the card reader channel with no end-of-job indication between jobs; RJS recognizes the beginning of each new job by the appearance of a JOB gard; For each job successfully appoind, the user will receive a confirming message on his densele. At the end of the stack, he must send an End-of-Data transaction to initiate processing of the last job; NETRJS will then close the channel (to evoid holding buffer space unnecessarily). At any time during the session, the user can recopen the cerd reader channel and transmit another job stack; He can also terminate the session and sign on later to get his output.

The user can abort the card reader channel at any time by closing the channel. NETRUS will then discard the last particily appoind job. If NETRUS finds an error it will abort the channel

by closing the connection prematurely, and also inform the user via his console that his Job was discarded (thus solving the race condition between EndmofmDate and aborting). The user needs to retransmit only the last Job, however, he could retransmit the entire stack (slithough it would be somewhat wasteful) since the CCN operating system enforces Job name uniqueness by immediately "flushing" Jobs with nemes already in the system.

If the user's process, NCP, or Host, or the Network itself fells during input, RJS will discard the job being transmitted. A message informing the user that this job was discarded will be generated and sent to him the next time he signs on. On the other hand, those jobs whose receipt has been acknowledged on the operator's console will not be effected by the failure, but will be executed by CCN.

The user may wait to set up a virtual printer (or punch) and open its channel until a STATUS message on his console indicates output is ready; or he may leave the output channel(s) open during the entire session, ready to receive output whenever it becomes evallable. He can also control which one of several evaluable lobs is to be returned by entering appropriate operator commends.

When RJS has output to send to a particular (virtual) terminal and a corresponding open output channel, it will send the output as a series of logical records. NETRJS will send an End-of-Data transaction and then glose an output channel at the end of the output for each complete batch job; the remote site must open a new channel to start output for another job. This gives the remote site a chance to allocate a new file for each job without breaking the output within a job.

If the user detects an error in the stream, he can issue a Backspace command from his console to repeat the last "page" of output, or a Restert command to repeat from the beginning of the lob, or he can abort the channel by closing his socket. If he aborts the channel, RJS will simulate a Backspace command, and when the user recopens the channel the job will begin transmission again from an earlier point in the same data set. If the user's process, NCP, or Host, or the Network (table fails during an output operation, RJS will set as if the channel had been aborted and the user signed off.

Although the NETRUS protocol was defined by UCLA CCN as an interim protocol for use to "get started" while the official RUE protocol was being designed, as noted in the previous section,

NETRJS has remained in operation and the official RJE has gone unimplemented. Further, any site wishing to compete with CCN for the batch service "market" is essentially obliged to implement NETRJS because the batch user processes use this (rether than the RJE) protocol. The lesson is quite sleer: (n an effortelimited world the first specification to be implemented has a very high probability of becoming the de facto standard.

2.5.5.3 Graphics

The eim of a graphical protogol is to allow users with various different kinds of display hardware at different sites in a network to make use of common graphids applications programs. One approach is a collection of special-spurpose protogols. Each application program could publish a description of the protogol needed to drive the program and to view the output. The prospective user would then write a program to interpret the published protogol and to drive his display (probably making use of existing graphics programming facilities at his site). This might permit a convenient division of labor between the gemputer executing the application program and the computer driving the display. The disadvantage of this approach is that the user must write a new program to interfece his display to each different

application protocol. In addition, there is no guarantee that the protocol required by the application program can actually be implemented within the user's operating system and display hardware.

Another approach is to develop one general-purpose protocol such as a Network Virtual Graphics Terminal that tries to provide facilities that a large number of application programs could use and that a large number of user sites could interpret as with TELNET and the NVT. Thus, one user program could be used to interface to a number of application programs. The disadventages of this approach are that the generality may preclude adequate response through the network, or that some application programs will find the general-suprose protocol too restrictive to be used at all. In addition, the design of such a protocol is not easy; attempting to provide a common device-independent framework for driving herdwere of qualitatively different capabilities may be very difficult.

The network graphics protocol attempts a middle course. It is not intended to satisfy all graphics needs, for all terminals, now and in the future. It is limited to ealligraphic pictures, to moderate interaction demands, and to "best effort" attempts to

generate the required graphics. Certainly, the impact of video technology will increase demand for variable character and maybe even enimetion techniques and it is anticipated that special epurpose protocols will be necessary for these (and gertain present) applications. In such cases, the general-purpose protogol can perhaps be used as a sterting point far development o f spedial-purposa protocols. The deneral-purpose protogol should be used whenever possible. however, so that separate user programs need not be written for each user site.

A user session with a graphics application program, as illustrated in Figure 2,5-11, is modeled on a TELNET session with a program; the user must log into a server system, execute several system commands, initiate the program, communicate with the program as it is running, perhaps interrupt a running program, log out, etc. A graphics session certainly requires all of these features; in addition it will require a pathway for transmitting graphics protocol information in both directions.

The application program, when it wishes to produce graphical output or to request graphical input, makes use of a server program (SP) which may simply be a subroutine package. The job

Figure 2,50111 Graphics Protocol

of the SP is to interface to the network graphics protogol on one side and to a "graphics language" on the other:

The protocol transmitted between the SP and the user program (UP) is the graphics protocol. It provides facilities so that:

- 1. The SP gan dause images to appear on the user's display sereen.
- 2. The UP cen report to the 8P any interactions that the user initiates, such as a key depressed on the keyboard or stylus interactions.
- 3. The SP can discover verious special properties of the user's display terminal, and can take adventage of them in conjunction with the application program.

These types of protogol trensmission are termed "output," "input," and "inquiry."

The job of the UP is to interpret the protocol and to generate appropriate devicemdependent information so that the image shown on the user's display corresponds to that specified by the SP using the protocol. It also implements the input and inquiry aspects of the protocol. The UP may have many "local

options" that are not covered by the protocol. For example, the UP might contain facilities for creating hard copies of the (mage on the display without engaging in any protocol exchanges.

Essentially the UP can be viewed as either an "interpreter of structured picture definitions" or an "interpreter of transformed picture definitions". In the protocol, the UP tolls the SP which kinds of formet it can implement. It is perfectly possible for the UP to implement both as the display images due to each of the formets are merged onto the screen.

These two different kinds of output format can be summarized as follows:

Trensformed: The protocol is used to build and modify a set of "segments" (sometimes dalled "regords") of a transformed display file, stored in the user Host. A segment is a list of graphical primitives that specify lines, dots, and text to be displayed at specify positions on the display screen; Individual segments may be deleted or replaced; they may be added to or removed from a list of segments to actually display on the screen (this is salled "posting" and "unposting" segments). If a picture is composed of many segments, changes to the picture can often be made by

replacing one or two segments; thus segmenting the display file helps to reduce the amount of information that must be transmitted through the network to effect change; Considerable experience with this type of picture definition has demonstrated that device independence can easily be achieved.

One advantage of transformed format is that the UP can be kept very simple; no transformations need to be performed in the User Host. A UP along the lines should be able to be implemented in an IMLAC. The burden of transformation is left to the server Host, presumably a large computer quite capable of being programmed to do transformations.

Structured: The protogol is used to build and modify "figures"; each figure is a sollection of "units". A unit may be a list of graphidal primitives, such as lines, dots and text; or it may specify a "eall" on another figure, together with a transformation to apply to the called figure. The protogol can replace individual units; altering a figure that is called in several places may cause widespread changes to the visible display. The structured format requires (in principle) even less network bandwidth for

updates then does the transformed format; many updates may involve changing a single transformation (e.g., to change the viewing transformation for a three-dimensional object).

Although a structured picture definition can be interpreted directly by a few display processors that have transformation ebility, most UPs that implement structured format will perform the transformation in software. This implementation is relatively difficults and certainly requires a fairly powerful user Host computer.

A remaining issue is when the displays generation software is to run in order to update the display. A useful technique is for the SP (and application program) to signal the UP whenever a splication (or batch) of changes is complete and the display should be updated. This technique has the following advantages:

- in order to update the displays never needlessly repeating transformations.
- 2. Screen erasures on storage tubes can be minimized; the screen will be erased a maximum of once per batch of updates, rather once per update.

All changes to the display appear "instantly". Network speed means that displayedile updates will arrive at the UP over a longish period of time. If the affect of these changes is delayed until a batch is finished, the display will appear to change "all at once", a much more satisfying effect than many slow changes. This is particularly true if an image is being replaced by another image that is a scaledwup version of the original.

To take full advantage of this technique, the application program should specify when the screen should be updated to represent precisely what is specified in the display file. These send patch updates dominands should precede each request for new user input so thus the user will see an upstordate image before formulating his response.

The problem of providing input facilities is even harder than that of providing output facilities. The difficulties are chiefly those of device independence and of adequate performance. The difficulty of achieving device independence is easily described; display hardware dan have a large number of very different kinds of input gadgets attached (light pens, tablet and

styll, Joysticks, knobs, buttons, stq.) that have different properties and different methods of reporting their output (e.g. periodically, on computer demand, or "whan something ghanges"). In addition, operating systems at user sites often enforce restrictions on the use of input aquipment in order to avoid undue system degradation. Further, if each input must be shipped to the server Hest, prosessed by the application program and SP, then any display updates shipped to the user Hest and processed by the UP before the user sees the response, it would be impossible to use many intersetive graphical techniques.

The device independence (saue is solved by inquiry) the UP reports to the SP a list of available devices. The SP and application program can then collaboratively arrive at an acceptable set needed for operating the application program. If the set of input devices is insufficient, the application program can perhaps engage in a dialog with the (human) user to seek remedies. Perhaps another version of the UP can be run which implements the required device or, if the application program and SP are sufficiently flexible, perhaps the commend language of the application program can be altered dynamically to permit its operation with the available devices. For example, if the UP responds that it has no coordinate input device (and they it is

not willing to simulate one with, say two knobs) then the AP might want to use a keyboardebased interaction sequence and to arganize the entire command system differently.

The performance diffigulties are addressed by permitting the SP to ask the UP to use some particular interactive technique in conjunction with an input device, and to report the results of the interaction. Such interaction techniques are termed "events". The techniques often involve providing "local feedback" so that the user sees the results of his interaction without a long network delay: Examples of local feedback ores displaying a "treaking dot" at the surrent losstion of a goordinate input device, or displaying a treil of "ink" behind the tracking dot, ats, Examples of the special "events" that the protocol provides are:

- Positioning; providing a pair of xey coordinates to apenify the position of something.
- Pointing; giving the coordinates of a displayed object in order to identify it (i.e., associate it with some other object or attribute).
- Strokings reporting a stream of coordinates identifying an arbitrary curve.

Draggings using a coordinate (or other) davice to cause some portion of the the display image to move in synchronism with the coordinate device.

protocol provides the SP with two basic methods for dealing with input devices: (1) to request and obtain the state an input davice, e.g., the current position of a soordinate input device, and (2) to enable various "events", end to obtain e report describing the "events" regulting from user actions.

The protocol has no set of stendard feetures; there is thus Network Virtuel Graphics Terminel as viewed by the protocol. The inquiry function is used to transmit to the SP a certain amount of detailed information about the terminal in use and the UP that drives it. This information will probably be requested by the SP when a session is initiated.

information transmitted by the inquiry response is in pert for information only. However, some of the information returned is essential in order for the SP to transmit legal protocol to the UP. In outline, the information returned is:

List of implemented protogol commends: This list tells the SP whether the UP implements transformed format, or structured format, or positioned text, or any combination of them, and so forth. In addition, this report tells which optional parts of the protocol ere implemented by the UP.

- Coordinate information: This information is necessary for the SP to darry out transformations that generate doordinates in the coordinate system used by the terminal.
- Parameters: These describe available character sizes, evailable intensity resolution, available line textures, etc.
- A list of available input devices and events: A "device number" is specified for each device; this is used when reading the state of a device, Similarly, an "event number" is specified for each event, and is dited when enabling or disabling it.
- An ASCII text strings This describes the terminal, e.g., "IMLAC PDS=1 in room 22".

The information in the inquiry response (that transmitted from UP to SP) that is not essential to further protocol operation may still be useful to the SP in order to drive the terminal intelligently. For example, inquiry can determine the kind of display being used, not so as to send device-specific

code to it, but so that the application program does not try to use a graphic technique on a terminal that sennot handle it (e.g., some sort of dynamics on a storage tube).

As may be obvious from the foragoing, the graphics protocol designers were unable to design a responsible "virtual" device; the variety in available hardware and software is too great, and no class of devices or techniques dominates the graphics field; Thus, as with FTP, there is a little of everything in the protocol. It seems clear that the development of a compact and useful graphics protocol is still very much a subject for research.

2.6 Performance (Incomplete section)

In this section, we discuss the many aspects of ARPANET performance.

2.6.1 Introduction

explanation of terms

short history

2.6.2 Theoretical Models, Measurements, Simulations

Kieinrocks* model

IMP measurement package/UCLA work

various experiments and measurement runs

2,6,3 Delay

In this section we consider what delay performance characteristics are possible in a given network. The topics discussed below include the identification of the components of delay, an analysis of the minimum delay possible for a round trip through a network, the reason for breaking messages up into packets, a brief look at queueing delay, and a discussion of delays for interactive traffic.

Components...of...Delay. We will first discuss the delay experienced by a single packet transmitted over a single hop. The components of delay which we will take to be fundamental variables are as follows:

1. Le propagation delay or speed of light latency.
This is the delay for the first bit of the packet to traverse the direct.

L s (direut length [mile])/(signal propagation rate [mile/sec]),

2. Te transmission delay.
This is the time for the bits of a packet to be clocked out on the circuit.

T = (number of bits in packet [bits])/(transmission rate [bits/sec]),

In the analysis below we will use Tp to denote T for a packet and Tr to denote T for a RFNM, which we will assume is a minimum length packet.

3. C m nodel processing delay.

This is the time it takes the node to process the pagket. It has two fixed components, dorresponding to the store operation and the forward operation, or receive and transmit, plus a random component due to queued tasks of higher priority. This component measures the interference experienced by packets queued for processor service. (In the following, we ignore delays at the source and destination nodes due to message processing.)

C = Cr[sec] + Ct[sec] + Cq[sec];
receive transmit queueing = rendom

A typical range of C (* probably 1 to 10 milliaeconds for an IMP=1ike node.

4. Do a queueing dalay

This is the time that the packet must wait for the transmission of the packets which precede it on the output queue, including

the output time of the packet currently being transmitted. Thus, Cq measures input queueing delay, waiting for central processor service, while Dq measures output queueing delay, waiting for dircuit service.

This is the time that the packet must wait in the event that its first transmission is unauccessful. This may happen if it was in arror or if the other node refused it for some reason, or, in the case of broadcast pircuits, there was a collision with another peaket.

In general, we will assume that the first two components are much greater than the last three. Tables 1 and 2 give some representative values for L and T_a (Note that a 50 Kbs direuit dan transmit 1000 bits in 20 ms, and that the first bit can travel about 4000 miles in that time, 80 a bit is about 4 miles long at 50 Kbs1)

Minimum Round Thin Delay. Now we can examine the minimum round trip delay, by taking the case of Cq # Dq = Dr = U_a Consider a message with P packate traveling over a path of H hops. If the delay at hop i is

Section 2 Design and Implementation

$$D(i) = L(i) + T(i) + C(i)$$

we can define the natural quantity D(ave), average hop delay, as follows:

$$D(ave) = (1/H) * zzzz(L(1)+T(1)+C(1))_{a}$$
 (1)

We can also define the less obvious variable D(max):

$$D(max) = xxxx(T(i)+Ct(i)), (2)$$

With these definitions, we can make the following two statements:

- The first packet experiences delay equal to H*D(ave).
- 2. The remaining P=1 packets follow through the network, each packet at most one hop behind the preceding packet, and these packets add (P=1)*D(max) to the total delay.

This enalysis can be illustrated by the following numerical example:

P # 3

H = 3

Ct = CP # 8.5

HOP

1 2 3

ARPANET Completion Report Chapter III

DRAFT Section 2 Design and Implementation

L Ø 1 3

1 2 3 2

Total delay # 21, as shown in Figure 2=35a,

Note that the bottleneck hop, in this case hop 2, has the largest T(i) + Ct(i), here equal to 3,5, but not the highest total delay. That is, hop 3 has total delay of 6, compared to the delay of 5 over hop 2. Note also that more than one packet is being transmitted at the same time, giving a pipelined effect, and reducing the total delay for the message.

This means that the delay for a single packet message is

$$D(SP) \neq \pi z z z z (L(1) + Tp(1) + C(1)), \qquad (3)$$

and the delay for a multipecket message is

$$D(MP) = (P=1) * Max[(=1,H](T(+)+C*(+)) + D(SP),$$
 (4)

The delay for a RFNM (s

$$D(RFNM) = zzzz(L(4)+Tr(1)+C(1))_{e}$$
 (5)

Therefore, the minimum round trip delay for a massage of P packets over H hops is

Figure 2=35 Example of Delays

$$D(MRT) = zzzz(2*L(1)+2*C(1)+Tp(1)+Tp(1))$$
 (6)

+ [P=1] *Max{i=1,H](Tp(i)+Ct(i)).

If the values of L, C, Ct, and T are the same for each hop, then we have a simplified minimum round trip delay $D(MRTS) \approx H*(L+Tp+C) + (P=1)*(Tp+Ct) + H*(L+Tr+C)_{1}$ (7) $\tilde{f}_{1}^{1} = t \text{ packet subsequent packets } RFNM$

Some survex are given in Figure 2=36 which illustrate the minimum round trip delay through a network for a range of message lengths and path lengths, for two sets of line speeds and lengths ever paths of 1 to 6 hops.

The Rationals for Rarketizing Messages. As a slight digression, we now consider the rationals for breaking up messages into packets in order to reduce delay. It will be shown that this rationals is similar to that for any pipelining technique, and the variables of interest will be identified. We will carry through this discussion for the simplified case of identical values of L, C, Ct, and T at each hop, first we rewrite equation 7 as

D(MRTS) = (P+1)*(Tp+Ct) + H*(2*(L+C)+Tp+Tr)

Figure 2=36 Min Round Trip Vs. Message Length

Note that the only one of the variables L, C, Ct, and T which is a function of packet length is T. Let us suppose that instead of a message of P packets, each of which has a transmission delay of T, the message is sent as a whole, with a transmission delay of PaT. The total delay can be calculated as above, donsidering this new entity as a long single-packet message. Using the equation for D(MRTS), we get an equation for minimum round trip delay, simplified, not packetized

$$D(MRTSNP) = H*(L+P*Tp+C) + H*(L+Tr+C)$$
 (9)

We can rewrite this as

$$D(MRTSNP) = (P=1)*H*TP + H*(2*(L+C)*Tp+Tr)$$
 (10)

Subtracting equation 8 from equation 10 we get a difference in delay

$$D(MRTDIF) = (P=1)*((H=1)*Tp=Ct)$$
 (11)

We make the following observations:

i. Clearly, if Pai, D(MRTDIF) # 0

There is no difference between the two techniques for single-packet messages.

2. Assuming Pair (f Hall D(MRTDIF) a = (P=1)+Ct < 0

That is, if HEL and PEL, it involves more dailay to break a message into packets because of the added processor overhead (generally small).

3. Assuming Tp>Ct, that the processor time is small, if P>1 and He1, then D(MRTDIF) > 0.

This is the usual case, and constitutes the basic reason to break messages into packets.

We can return to the numerical exemple used earlier with the equivalent conditions for a single packet message as long as the 3 packet message considered above:

P = 1

H = 3

Ct # Cr # 0,5

HOP

1 2 3 total

C 1 1 1 3

L 0 1 3 4

T 6 9 6 21

total 7 11 10 28

Total delay # 28, as opposed to 21 for the packetized case (see Figure 2=35b).

Queueing Delay. We have examined the minimum round trip delay as a function of message length, network path length, and the packatizing strategy. It is appropriate at this point to analyze the effects of additional delays which may be present. To perform the minimum delay analysis, we made 3 assumptions, each of which should be remaxamined at this point:

- i. Cq = 0. The packet arriving may have to wait on an input queue before it is serviced by the processor. This time is generally quite small, but it is a rendom variable which may take on large values if there are time-consuming higher-fority tasks in the system.
- In the analysis of actual delays in a network. A great deal of theoretical work has been done in studying queueing delay, particularly by Kleinropk, who has used both analysis and simulation in this regard. For the simplified discussion at hand, it is sufficient to note that each packet on the queue adds an additional T+Ct to the delay of all packets behind it on the queue.

In the type of network and line being considered. For some wideband circuits, particularly setellite channels, the error rates are very low and few retransmissions may be necessary for reasons of errors. However, satellite links may be used in a broadcast competition mode so that some packets are lost in collision with others, finally, there is always some chance that the adjacent node will refuse to accept a packet for lack of processing resources.

Low Delay for Interactive Traffig. Perhaps the most important consideration about delay in a network is this some traffic consists of interactive, high-priority messages and this traffic must be delivered to its destination as rapidly as possible. This is in contrast to bulk trensfer traffic which is not so delay-critical. The most obvious case of sugh interactive traffic is most man-computer dialogue, which consists of rather short messages between the computer and a man at a terminal, Here there is a definite threshold for delay, Below this threshold, delay is acceptable, and above it dalav unacceptable. There is no added benefit if the delay is gonsiderably below the threshold, and it is likely that once the delay is much above the threshold, almost any value of delay (5 equally unacceptable.

DRAFT

Section 2 Design and Implementation

Another way of looking at the bimodal nature of network traffic is to consider that much of the delay for an interactive message is in the network itself. That is, it is generated and quickly sent into the network. When it is delivered, it will be processed quickly at the destination Host. Bulk transfers on the other hand may experience lengthy delays outside the network due to buffering considerations and the very size of the data (secondary storage or tapes or cards may be involved in the data transfer, greetly ingressing delay).

2.6.4 Throughput

In this section we consider what throughput performance characteristics are possible in a given network. The topics discussed below include an analysis of nodel processor bandwidth, an examination of circuit overhead, a quantification of the buffering in peckets required for a network line and the buffering in messages required for a network path, the throughput requirements of bulk transfer traffic, and the tradeoff between delay and throughput.

The Firetive Bandwidth of The Node Brogesor. We will first examine the effective processing capability of the node somputer in a network with veriable message length. Some numerical examples will be given to support the intuitive notion that the processor is most effective for long messages, given some very general assumptions about the packet processing involved.

We begin by defining some new quantities of interest in studying throughput in networks, The quantities which we will take as fundamental variables are as follows:

1. Bd = the number of host data bits in a packet,

- 2. 8s = the number of softwere overhead bits per packet. These bits include header information such as address, and identifying information such as message number.
- 3. Bh = the number of hardware overhead bits per packet. These are typically framing bits for the circuit, and error detection bits such as checksums or redundant information.
- 4. P = the number of packets per message, as above.
- 5. Boot # the total number of host data bits per message.

Btot & PeBd = unused bits in the lest packet,

Now we will examine how long it takes the node processor to store and forward a message. As noted above, we ignore the processing at the source node and at the destination node. The time required to process a storemend-forward message is a function of the following parameters:

6. C s the packet processing time, as above,

We assume that C is independent of packet length or type. To the extent that this is not true, the length-dependent component of C can be accounted for in item 8 below.

 7_s 8WPo = the fraction of the bandwidth of the processor taken by overhead.

Due to certain necessary periodic processes within the node, notably the routing computation, effective processor bandwidth is reduced.

8. BWIO m the I/O rate of the node in bits/sec.

We assume here that BWIO is a linear function of the number of bits in the message. In most I/O architectures, it is probably a function of the number of computer words in the message, which is identical apart from unused bits in the last word. We are also essuming that the I/O transfer steels cycles from the processor, reducing its effective bandwidth.

9. I m the I/O transfer time in seconds. We will denote the I/O time for a packet, a message, and a RFNM as Ip, Im, and Ir respectively:

 $Ip = 2*{8d+8m}/8WIO$

Im = 2*(Btot+P*Bs)/BWID

IF = 2+8s/8HIO

Saction 2 Design and Implementation

10, MT m the total time taken to process a message.

 $MT = C \pm (P+1) \pm (1 + BWPo) + Im + Ir.$

11. BWPd = the maximum data bandwidth that the node can support.
BWPd = Btot/MT.

This the the number of host data bits per second that the node can process.

12. 8MP) a the maximum | the bandwidth that the node can support.

 $BWP1 = \{Btot + (P+1) + (Bs+Bh)\}/MT$

This represents a processing dapability limit on the number and speed of the circuits that can be connected to a node, The difference between the two quantities, BWP1=BWPd, is a measure of the line overhead at a given message length.

At this point, a numerical example may be illustrative. In the ARPANET the programor bandwidth of the IMP, both BWPd and BWP1, can be plotted as functions of the message length, and some results are given in Figure 2=37. (The discontinuities indicated in the curves are the result of packet(zing messages.)

Figure 2=37 Processor Bendwidth Vs; Message Langth

The Effective Bandwidth of the Network Circuits. In this section we will consider some of the factors acting as overhead to reduce the effective circuit bandwidth for network lines. That is, we wish to catalog all the kinds of transmission that take place on the network lines that are not satual host data bits, and from this accounting determine which factors are the dominant ones. There are three basic kinds of line overheads

1. Line overhead in bits/packet,

We have detailed two components of this overhead earlier, Bs and Bh.

2. Line overhead in bits/message.

In the terminology we have been using, the RFNM is overhead on a message basis, and under our assumption that it is a minimum length packet, it contributes a number of bits equal to Ba+Bh.

3. Line overhead in bits/second of network system traffic.

Some examples of this kind of traffic can be given, along with the approximate order of magnitude for the traffic rate in the ARPANET:

routing 1800 bits/sec

line alive/deed 100 bits/sec

NCC status reports 10 bits/sec

core rejoads and dumps 1=10 bits/sec

We will make the assumption that routing messages are the primary source of traffic other than messages between hosts. For this reason, we will define the variables

Br = the number of bits in a routing message

Fr m the frequency of routing messages [[/sec]

BWCs = Br#Fr = the bandwidth of the girquit

and ignore the other components of periodic everhead,

We can now total the overhead from all considerations, by converting to a common dimension, bits/sec. In order to do this, we must introduce another variable,

Fp m the number of packats per second.

We then son express the total bandwidth of the sirguit given to overhead as

8WCo m Br*Fr + (Ba+8h)*Fp + (Ba+8h)*Fp/P routing packets RFNMs

We can rewrite this as

BWCo = BWCr + (Bs+Bh)*Fp*(P+1)/P

In comparison with this overhead rate is the agruel data rate on the dircuit,

BWCd = Bd*Fp

We can ovaluate the fractional overhead percentage as

BWCo/BWCd = (BWCe)/(Bd*Pp) + (Bs+Bh)*(P+1)/(Bd*P)

Another quantity of interest is the maximum data rate that can be attained for given values of the systems parameters. Given a circuit with bandwidth BWC, the maximum data rate occurs when Fp is at a maximum,

FP(Mex) = (BWC+BWCe)/Bd

Substituting the expression for SWCo above, we get

FP(Max) = (BWC + Be + FP)/(Bd + (Bs + Bh) + (P+1)/P)

The numerator indigates that the routing message bandwidth comes off the top, leaving a reduced effective bandwidth, which is then

used for both data bits and packet and massage overhead bits:
The graphs in Figure 2-38 show the behavior of some of these variables; the variable plotted is Fp(Max)+8d, that is 8WC-8WCo, which is the maximum data rate for a given massage size.

The Recket Buffering Required for Network Circuits. We now turn to an examination of the number of packet buffers required to keep a communications circuit fully loaded. This number is a function not only of line bandwidth and distance but also of packet length, nodel delays, and acknowledgment strategy. We will assume that the node buffers each packet that it transmits until it receives an acknowledgment, meanwhile transmitting other packets to utilize the circuit efficiently. If it does not receive the acknowledgment in the expected time, it retransmits the packet, We also assume that packet buffers are a fixed size large enough to hold the largest packet. The expected time for an acknowledgment to return is the sum of:

- Is To π the transmission time for the packet, a function of line bandwidth.
- 2, i = speed=of=light delay for the first bit of the packet to errive at the other node, a function of line length.

Section 2 Design and Implementation

Figure 2+38

Section 2 Design and Implementation

- 5. C = Cr+Ct = the processing delay in the other IMP, to receive the packet and return the acknowledgment.
- 4_{\pm} Dq = queueing delay for the returning asknowledgment, the time it waits for any other transmissions sheed of it.
- 5. Ta # the transmission time for the acknowledgment,
- 6. Le speed=of=light delay for the first bit of the acknowledgment to arrive at the first node.
- 7. Cr = the processing delay for the acknowledgment,

Our first simplifying assumption is that the processing times are small compared to the other delays and can therefore be ignored:

We can then state that the minimum number of peaket buffers needed to keep a girquit fully loaded is

BFp =
$$(Tp + L + Dq + Te + L)/Tp$$

This can be rewritten in somewhat more meaningful form as

$$BFp = 1 + 2 \times L/Tp + (Dq + Ta)/Tp$$

This expression indicates that one buffer is always necessary, to account for the packet transmission time itself. More buffers may be required if the direct is long compared to the packet transmission time, or if the adknowledgment transmission takes a long time compared to the packet. Stated differently, the number of buffers needed to keep a line full is proportional to the length of the line and its speed, and inversely proportional to the packet size, with the addition of a constant term.

In order to proceed further with the enalysis, we need to introduce two new terms:

Ta m the transmission time of the shortest allowable packet

It is the transmission time of the longest allowable packet

We also need to postulate a traffic mix of long and short packets, with x/y the ratio of short packets to long packets in the channel. Now we can define Dq and Ta in terms of Ta and Tl: We make a worst-case assumption for Dq:

Do a Ti, the acknowledgment has highest priority (equivalently, it piggybacks on all packets), but it must wait for the transmission of a maximum-length packet which has just begun.

The assumption for Ta is rether arbitrary;

Tam(Ts+T1)/2, the acknowledgment piggybacks on an Anna Parage Tamping packet.

We now state the result for the number of packet buffers required given the above set of assumptions:

$$BFp = 1 + (2*L + T1 + (T*+T1)/2)/Tp$$

Using the ARPANET values of Ts and T) given above (n Table 2, and shoosing a variety of line lengths and traffic mixes (shown as the ratio of short packets, 5, to long packets, L), we can present some numerical results as a family of survex shown in Figure 2=39. Note that the knee of the survex occurs at progressively shorter distances with ingressing line speeds. In fact, if we define the knee to occur when the linear term is equal to half of the constant term, then the knee occurs when

$$L = (Ts + 2*Tp + 3*T1)/4,$$

or for a line length of 225*10**6/8WC miles. The constant term dominates the 9.6 Kbs case, and it is almost insignificant for the 1.4 Mbs case. Note also that the separation between members of each family of surves remains constant on the log scale, indicating greatly increased variations with distance.

DRAFT

Section 2 Design and Implementation

Section 2 Design and Implementation

Figure 2=39

Ither Massage, Buffering, Required, for Network Raths. This section takes up a topic which closely parallels that of the last section. Here we will examine the number of massages needed to obtain full bandwidth over a network path of many lines. That is, we will compute how many massages must be in flight between two nodes in order to keep all the intermediate lines fully loaded. Actually, the best one can do is to keep all the lines in the path of the lowest direct bandwidth fully loaded. It turns out that this analysis is quite simple given all the definitions of the preceding sections. The number of message buffers needed is computed by taking the round trip delay for a message and dividing it by the time taken to transmit a single message. That is,

BFm = D(MRT)/(P+(To+Ct))

Uning equation 6 in section 2,2,1, we can obtain a more detailed axpression for the simple case of equal delay at each hops

 $8Fm = \{P=1\}/P + H*(2*(L+C)+Tp+TP)/(F*(Tp+Ct))$

It is clear from this expression, and on intuitive grounds, that SFm is a minimum for maximum length messages, that is for large P. The curves in Figure 2040 show the dependence of SFm on line characteristics and on the length of the network path, for PES

DRAFT Section 2 Design and Implementation

ORAFT Section 2 Design and Implementation

Figure 2-40

and ARPANET values of the parameters. For each of four line speeds, the buffering requirements are plotted for network paths made up of a number of land lines (the length of the lines (squen with each surve). Also shown are the requirements for the same network paths with the addition of one satellite link running at the same bendwidth as the land lines.

The consequences of the above discussion are several. For the communications subnetwork, it means that the nodes must be able to do bookkeeping on several messages in flight between two nodes. Further, the network must buffer these messages in the memory of the source or destination node for the duration of their flight, in addition to the packet buffering that takes place instantaneously at the intermediate nodes along the path, This has some important remifications for the design of the software for the node somputer. A second say of issues is the affect of message buffering on host computers, It is given that the communication protocols that the hosts use must also be engineered to the parameters of the network, if they are to obtain full throughput levels.

High.Throughput.for.Lang.Data:Transfers. The several topics examined in this section all point to a single conclusions the larger the packets in a message, and the larger the messages in a

data transfer, the higher the level of throughput that is potentially attainable. For reasons of processor overhead, ejecuit overhead, and buffering considerations within the nodes, it is always better to have long packets and messages if high data rates are desired.

The Tradeoff between Low Delay and High Throughout. It is alear that the two goals of low delay and high throughput are often in conflict because the node has limited resources with which to service its hosts and lines. It is difficult to quarantee both low delay and high throughput to several competing sources. The approach taken in the ARPANET is as follows: the IMP program has been designed to perform well under bimodal conditions. It provides quick delivery for short traffic interactive measures and high throughput retes for long files of data. This optimization of the program for a specific model of traffic behavior occurs at many levels, and is essential for belanced performance characteristics. (Other conflicts exist, such as the conflict between high throughput and congestion/deadlock control mechanisms.)

2,6,5,1 IMP Techniques

The network is designed to be largely invulnerable to directly or IMP failure as well as to outages for maintenance. Special status and test procedures are employed to help gops with various failures. In the normal dourse of events the IMP program transmits hellos (routing messages). The adknowledgement for a hello packet is an Isheard-you (IHY) bit in a returning null packet.

A dead line is detected by the systemed absence (approximately 3,2 seq) of IHY messages on that line, No regular packets will be routed onto a dead line, and any packets awaiting transmission will be rerouted. Routing tables in the network erredjusted automatically to reflect the loss. Receipt of consecutive I=heardeyou packets for about 30 seconds is required before a dead line is defined to be alive once again. The IMP program takes into account the fact that an IMP may have only one working line to the network. In this case, a line may make twice as many errors as is usually permitted before it is declared unusable. This mechanism is an attempt to improve natwork availability from singly-connected sites.

A dead line may reflect trouble either in the communication facilities or in the neighboring IMP (tself, Normal line errors daused by dropouts, impulse noise, or other similar conditions should not result in a dead line, because such errors typically last only a few milliseconds, and only occasionally as long as a few tenths of a second. Therefore, it is expected that a line will be defined as dead only when serious trouble conditions occurs.

If dead lines eliminate all routes between two IMPs, the IMPs are said to be disconnected and each of these IMPs will discard massages destined for the other. Disconnected IMPs cannot be repidly detected from the delay estimates that errive from neighboring IMPs. Consequently, additional information is trensmitted between neighboring IMPs to help detect this sondition. Each IMP transmits to its neighbors the length of the shortest existing path (i.s., number of IMPs) from itself to each destination. To the amellest such reselved number 2007 destination, the IMP adds one. This incremented number is the length of the shortest path from that IMP to the destination. If the length ever exceeds the number of network nodes, destination IMP is assumed to be unreachable and therefore disconnected.

Messages intended for deed hosts (which are not the same as dead IMPs) cannot be delivered; therefore, these messages require special handling to evoid indefinite circulation in the network and spurious arrival at a later time. Such messages are purged from the network at the destination IMPs. A host computer is notified about another dead host only when attempting to send a message to that host.

The components of the IMP program dedigated to improving reliability have two main functions. First, the software is built to be as invulnerable as is possible in practice to herdware failures. Second, the software isolates and reports what failures it can detect to the NCC. With intermittent failures, it is important in practice to keep the IMP program running and diagnosing the problem rather than keeping the IMP down for long periods to run special herdware diagnostics.

The IMPs use the technique of softwere checksums on all transmissions to detect errors in paskets, protecting the integrity of the data and isolating hardware failures. The end-to-end softwere checksum on packets, without any time gaps, works as follows (refer to Figure 2-41).

Flaure 2=41 Softwere Checksums

- e= A shecksum is computed at the source IMP for each packet as it is received from the source host (interface 1).
- The checksum is verified at each intermediate IMP as it is received over the circuit from the previous IMP (interfaces 3 and 5).
- -- If the checksum is in error, the packet is discerded, and the previous IMP retransmits the packet when it does not receive an acknowledgement (interfaces 2 and 4).
- The previous IMP does not verify the checksum before the original transmission, to gut the number of checks in helf. But when it must retransmit a packet it does verify the checksum. If it finds an error, it has detected an intrasIMP fellure, and the packet is lost. If not, then the first transmission was lost due to an intersIMP fellure, a circuit error, or was simply refused by the adjacent IMP. The previous IMP holds a good copy of the packet, which it then retransmits (interfaces 2 and 4).
- After the pecket has suggestfully traversed several intermediate IMPs, it arrives at the destination IMP.

 The checksum is verified just before the packet is sent to the host (interfece 6).

This technique provides a checksum from the source IMP to the destination IMP on each packet, with no gaps in time when the packet is unchecked. Further, the length of each packet is verified. Any errors are reported to the NCC in full, with a copy of the packet in question. This method enswers both requirements stated above; it makes the IMPs more reliable and fault etalerant, and it provides a maximum of diagnostic information for use in fault isolation.

One of the major questions about such approaches is their efficiency. We have been able to include the software checksum on all peckets without greatly increasing the processing overhead in the IMP's. The method described above involves one checksum calculation at each IMP through which a packet travels, we developed a very fast checksum tachnique, which takes only one instruction execution per word. The program computes the number of words in a packet and then jumps to the appropriate entry in a chain of instructions. This produces a simple sum of the words in the packet, to which the number of words in the packet is added to detect missing or extra words of zero, with the inclusion of this code, the effective processor bendwidth of a Honeywell IMP is reduced by one-eighth for full-length store-end-forward packets. This add checksum is not a very good

one in terms of its error-detecting capabilities, but it is as much as the IMP can afford to do in activare. Furthermore, the

primary goal of this modification is to assist in the remote

diagnosis of intermittent hardware failures,

A different set of reliability measures has been instituted for routing. It is clear that datestrophic effects den follow for the network as a whole when a single IMP begins to propagate incorrect routing information. This failure may be due to a memory failure in the data area or in the program itself. A single broken instruction in the part of the IMP program that builds the routing message causes the routing messages from the IMP to be rendom data. The neighboring IMPs interpret these messages as routing update information, and traffic flow through the network can be completely disrupted.

It is useful to provide a brief look at some of the problems endountered in the ARPANET with the reliability of routing. In 1971, there was an IMP which drashed every few days in such a way as to gause all the other IMPs in the network to majfunction as well. It was finally determined that its core memory was faulty and sometimes the routing messages read out from memory by the modem output interfeces were all zeroes. The adjacent IMPs interpreted such an erroneous message as stating

that the IMP in question had zero deley to all destinations = that it was the best route to everywhere; Once this information was propagated to the other IMPs, the whole network was operating with incorrect routing information, and normal operations came to a halt. The problem here is that the information that one node intended to transmit was not suggestfully received by another node.

A different problem happened in 1973, with similar symptoms of the whole network being affected. It was traced to an IMP having an incorrect instruction, due to a memory failure, in the middle of its routing program. This had the effect that the IMP computed an incorrect pointer to where it thought the input routing data was in memory, and so it proceeded with the routing computation on the basis of random data, and then propagated this information to the adjacent IMPs. This meant that the routing throughout the network was continuously changing, oscillating wildly, as the failed IMP sent out the incorrect routing data. Here the problem was that the node was not running the correct routing process, but it too had global effects.

A finel problem happened when one IMP, at Aberdeen, suffered a dropped bit in memory in its route directory. The

entry in the directory for the effected destination IMP, UCLA (many hops and 3000 miles away), was then zero. A zero entry in the directory at an IMP means that the entry is for the IMP itself. Thus, the Aberdeen IMP took all traffic for the UCLA IMP as traffic for itself. The dropped bit in the Aberdeen IMP meant that part of the network could not communicate with the UCLA IMP, though there were no failures in that part of the network or in the UCLA IMP. Note that the route directory is a completely internal table; it is not exchanged with the adjagent IMPs, and yet a failure in this table had global effects.

There is one lesson to be drawn from such ingidents:

The routing algorithm, is wiremely, indectant, to metwork reliability....sines...it...it..meliumations...the maked the routing algorithm...the property...thet...Teli...the...nodes...must...be...performing...the computations are forming...the comput

Having described seme of the problems that can arise in the course of equal network operations, we now turn to a more systematic evaluation of the problem areas and possible solution techniques.

In attempting to catalog the kinds of errors that can be introduced into routing data, it is useful to recell the system components that may cause the errors, and how these domponents fail. The processor, memory, I/O interfages, and the circuit between nodes are all hardware components that may have either a solid or intermittent failure, and the routing program or the routing data may be affected. In transferring a set of data from one node to an adjagent node, there are the following possibilities for errors

- is Incorrect data stored in memory at the source: The program may store it in the wrong locations in memory; the memory may fell at the data locations:
- 2. Date changed between dreation and transmission: The program may mistakenly change some of the date; due to hardware or software failure: the memory may fail:
- 3. Incorrect transfer of data from memory to line interface at the source: The program may direct the hardware to send the data from the wrong logations; the hardware may make the data transfer with errors.
- 4. Incorrect transmission of the data over the directly The circuit may be noisy, and make data errors; the directly may be unusable.

- 5. Data changed between reception and use: The program may mistakenty change some of the data, due to hardware or software failure: the memory may fail.
- 6. Incorrect transfer of data from the line interface to memory at the destination: The program may direct the hardware to read the data into the wrong locations; the hardware may make the data transfer with errors.
- 7. Incorrect data read from memory at the destinations. The program may read the wrong locations in memory; the memory may fail at the data locations.

We now explain briefly the measures taken in the ARPANET to protect against the failures cataloged above. The central aims of the approach taken in the ARPANET are:

- 1. To detect errors of any kind that may affect the routing computation.
- 2. To perform error checking often enough and in enough places to provide fault isolation in the event of failure.
- 3. To report as much data as possible about the failure to the Network Control Center for later diagnosis.
- 4. To allow the node which detects a failure to continue to function, whenever possible. This includes taking remedial action as appropriate.

5. Most important, to localize the effects of a failure by ensuring that incorrect routing information is never exchanged between nodes.

The basic technique used in the ARPANET to improve the reliability of the IMPs and enhance their ability to provide diagnostic information is checksumming. The original ARPANET design included checksum hardware in the modem interfaces, to detect errors in packets due to noisy direction. This concept has been extended greatly as follows:

- Is Routing messages sarry a special checksum (computed in software in the present IMP, but there will be special hardware for this purpose in the next generation IMP) which is distinct from the checksum on transmissions over the circuits.
 - The process which constructs each routing message builds the checksum. The sheeksum refers to the intended contents of the message, not the actual contents. That is, the program which generates the routing message builds its own checksum as it proceeds, not by reading what has been stored in the routing message area, but by forming the checksum of the intended contents for each entry as it computes them. This must be done in software, though the next two can be done in hardware.

- b. The process which sends out routing messages always varifies that the checksum is correct before it transmits a message. This detects intrasnade failures such as memory failures, processor failures, and random changes made to the data.
- o. The process which eccepts routing messages always verifies that the chacksum is correct before marking it as acceptable data. This detects inter-node failures such as memory transfer problems, interface noise, and so on.
- 2. All routing programs are chacksummed before every execution, to verify that the code about to be run is correct. The checksum of the program includes the preliminary checksum computation itself, the routing program, any constants referenced, and anything else which could affect its successful execution. These programs include:
 - a. The process which constructs routing messages,
 - b. The process which transmits routing messages,
 - o. The process which accepts routing messages,
- d. Any periodic processes which affect the routing data.

 In each case, if an error is detected, the program is immediately reloaded from an adjacent node, and a diagnostic message is sent to the NCC indicating the kind of failure.

3. Experience in Use and Operation

The actual network, the objectives and design of which have been discussed in such detail in the preceding sections, has changed during its lifetime from a computer communications research vehicle to an operational system supporting a wide variety of diverse users, pursuing independent computational goals in their deily work. This section briefly discusses some of the communities of interest which have been involved in ARPANET utilization, some of the novel systems which have been built to take advantage of network fasilities, and some of the ongoing communications experiments which take the existence of the ARPANET as a base upon which to build. The section concludes with a discussion of the operation of the network, since the network operation and maintenance is heavily relient upon the network itself.

A discussion of the Host experience with the ARPANET should surely begin with mention of some of the first natwork uses. Undoubtedly the very first use of the network for a task which was not some aspect of monitoring the network itself was the conversion of the Stanford Research Institute (SRI) On-Line System (NLS) software from an XDS-940 to a PDP-10. Conversion

began several months in advance of the delivery of the DEC PDS=10 system to SRI with the modification of the XD54940 compiler to produce PDP=10 object code. By early 1970 this modification had been performed and object code was ready for sheckout. initial checkout was carried out by sending the object code through the ARPANET to the PDP=10 system at the University of Utah, using a precursor of the standard Hostwidgst protocol. The code was executed at Utah and a record of its execution history (trades, console outputs, etc.) stored in the file system for later transmission back to SRI, This network use was important both because it showed on early, significant, payoff from the investment in the network, and because it demonstrated the feasibility of communication between dissimilar Host computers. In fact, it should be noted that there was a PDP=10 computer installed elsewhere within SRI at the time this network use was taking place, but it proved to be much more convenient to use the Utah PDP=10 because both it and the SRI 940 were connected to the ARPANET, while the SRI PDP=10 was not.

Several early "demonstration project" uses of the ARPANET were documented in midw1970 [Cerr70] and early 1971 [Metes]fe71];
These included a number of different examples of terminals directly connected to one Host being used to access a

time=sharing system on some other Host (e.g. SRI XDS=940 to Utah PDP=10, UCLA Sigma=7 to SRI 940, PDP=6/10 to Multics at MIT, and MIT PDP=6/10 to Hervard PDP=10). A more embitious experiment involved sending graphics programs and 3D data from Harvard's PDP=10 to the MIT Dynamic Modeling/Computer Graphics complex which included an Evens & Sutherland Line Drawing System, a quite fancy graphics processor. The data was processed repeatedly on the E&S system and 2D scope data was returned to Harvard's PDP=1 for display, This setup was used to run a simulated landing of an airplane on an aircraft carrier under control of user console input data. The importance of these experiments was that they demonstrated the ability of the network to support interactive traffic without excessive delays, although the aircraft carrier simulation was bendwidth limited. It was noted in [Metcalfe71], in a description of the Harvard/MIT console gonnection, that:

"...it was found that the response from the Harvard system at MIT=DM/CG was seamingly as fast as could be expected from one of their own donacles. This fact is particularly exciting to those who don't have a feel for network transit times when it is pointed out that such response was generated through two

time=sharing systems, three user level processes, and three IMPs, all connected in series."

Experiments such as these provided the impetus needed to hasten the development of protocols (and NCPs), while also providing valuable tests of some of the original protocol concepts and of the subnetwork of IMPs.

Since the time of these experiments network use has increased tremendously. The following sections briefly describe some of the details of this use. It should be noted that these sections draw heavily on the facts and concepts presented in the cited references, and in some cases incorporate text as published in the references without explicit indication.

3.1 Server Systems

It has become common when discussing the ARPANET Hosts to categorize them as either "servers" or "users" depending, obviously, upon whether they tend to provide services to other sites or to provide access to the network so that local people can access remote services. Of course, many server systems also provide an access path for local users to reach the ARPANET, so that there are almost no pure server systems (e.g., the Datacomputer and the ILLIAC IV), Similarly, there are few user systems which offer no services, other than the TIPs. However, most systems can be fairly unambiguously datagorized as servers or users.

It is worth noting that size alone does not determine whether a system is a server or a user. For example, the PDP=11 based Saismic Input Processor is a pure server system which acts as buffer between continuous low speed inputs and occasional high speed transfers to the Datagomputer. On the other hand, the CDC 5500 at Fleet Numerical Weather Central is an ARPANET user system which is used only to gain access to remote network services.

Among the network server systems there is quite evident functional specialization, with surprisingly little competition within each specialty. Thus, it is appropriate to think of the network servers as a complementary set of resources, among which choices are made on the basis of the Job to be performed. The primary exception to this situation is in the area of general purpose time-sharing with an emphasis on text handling. In this field there is some competition between the various PDP=10 systems on one hand and the Multics systems on the other; there is also competition among the verious PDP=10 centers, notably the TENEX centers at BBN and ISI.

Some of the natwork server systems were in existence (or, at least, were coming into existence) prior to the construction of the ARPANET and are "incidentally" connected to it; such systems include the UCLA Campus Computing Network (CCN) IBM 360 Model 91, the MIT Multips system, and the BBN Research Computer Center (RCC) TENEX systems. Others, such as the Datadomputer, included ARPANET connection as a relatively important part of their planning. (In addition, although its development was started well before the ARPANET, the planning for access to the ILLIAC IV dame to rely increasingly on the ARPANET,) Thus, some servers were relatively indifferent to whether they attained significant

outside usage or not, some had a readywarde "captive" network user community, and some actively sought new users. Further, some of the centers were run on a strictly economic basis, while others were as much (or more) influenced by research goals as by economics. For all these reasons, the server organizations interfaces to the users varied widely, from eager helpfulness through relative indifference.

Regardless of the attitude of the server organizations, however, none of them had significant prior experience dealing with a user community spread over a very wide geographic area; The response modes which had been developed to provide user services on a campus, or within a few office buildings of a single company, or within a single city and its suburbs proved inappropriate to a user community spread across the entire continent (not to mention England, Norway, and Haweii). In particular, the mode of herdcopy documentation which assumes a computer center builstin board where "recent changes" are posted, a donaultent's office nearby, and other system users in nearby offices, proved to be unworkable.

One of the first issues to be addressed by the ARPANET community was the question of how to discover what programs and

databases existed at what servers. An "ARPANET Resource Handbook" was generated at BBN and later taken over by the Network Information Center (NIC) at SRI. This handbook is loosely enalogous to the "yellow pages" published by the telephone company, providing a single place where the providers of service can announce what is available. The Resource Handbook is a logical companion to the ARPANET Directory, also produced by the Network Information Center, which is analogous to the telephone company's "white pages", a listing of how to contact individuels by address, telephone, and (more importantly) vis "network mail" (see Section 3.4).

The Resource Handbook, however, only provides pointers to existing resources; it doesn't provide much help in learning to use them, nor does it help in case of difficulty. The server systems have generally adopted similar strategies for dealing with these continuing usermessistance issues; strategies which are extensions of ideas which originated in logal timeshering systems. They (notude the concepts of "linking" two terminals (e.g., a user and an operator), "mailboxes" where requests for help can be left and answers collected, "system messages" sent to the console of each user, and oneline documentation.

In addition, some common strategies have been developed for the deses when human interactions are desirable. At the urging of the NIC, each site nominated a "Technical Lieison", competent to provide answers to system questions to a fairly detailed level, whose name was publicized in both the Resource Handbook and the ARPANET Directory. Most server sites also expanded their ability to provide consulting vie telephone; some sites, including the Network Information Center, provided tollarine voice telephone service for user consulting. Some of the servers, notably CCN, also encouraged user visits to the server site at least once, and provided periodic staff visits to the larger users [Kehi73].

3.1.1 UCLA = Campus Computing Network

The Compute Computing Network (CCN) of UCLA has been operating a large server system on the ARPANET since mid=1971; CCN operates on IBM 360/91 CPU with 4 megabytes of high-speed memory and a large secondary-storage configuration. For a number of years, this was the largest and fastest general-purpose computing system on the ARPANET.

As an ARPANET server, the 360/91 has been primarily used for large-scale numerical computation or "number crunching", in a

batch-processing mode. The CCN system also provides TELNET access to the IBM general-purpose timesharing system TSO and implements the file Transfer Protocoly both of these have been most useful as adjuncts to batch processing. In addition to its computational power, the CCN system offers users an extensive application program library encompassing a wide variety of fields, many of which are unavailable elsewhere on the ARPANET.

CCN's chance to obtain a connection to the ARPANET was a result of the presence at UCLA of Professor L. Kleinrook and his students, including 8. Crooker, J. Postel, and V. Cerf. This group was not only involved in the original design of the network and the Host protocols, but also was to operate the Network Measurement Center (NMC). For these reasons the first delivered IMP was installed at UCLA, and ARPA was thus able to easily offer CCN the opportunity for connection.

CCN reportedly had several interests in becoming an ARPANET Server System, including:

In The apportunity to become involved in an important new aspect of computing technology, helping CCN to obtain a lead in technical expertise.

- The prespect of additional income, coming at a time when the recession was causing cutbacks of federal and state funding, but an early point in the 360/91 life cycle with only about 50% loading generated locally. Further ARPANET users were expected to need betch number-grunching, an ideal load for CCN because of its characteristics of being CPU-bound and suitable for off-shift scheduling.
- The expedted local Benefit to students and faculty of having new application packages and programs installed on the 360 which would not otherwise be available.
- 4. The benefits to local users of access to remote resources, both computational and human, especially from the point of view of becoming linked into a major research community.
- 5. The intengible but real benefit of the prestige according to sites associated with the ARPANET.
- Of course, many of these reasons pertain to other server systems as well.

Although ARPA provided the UCLA IMP and paid for the necessary CCN interface hardware in order to encourage the availability of number-crunching services CCN incurred

significant costs in developing the necessary Host software. These costs could only be recovered by attracting new users, Thus, CCN was especially interested in seeing that the ARPANET protocols developed to support the types of service CCN affered were officient, reliable, and easy to implement as well as being specified early. This frequently led to conflicts with groups involved in the design of protogols, who may have exhibited more of a desire for elegance and generality, even at the cost of more difficult implementation and less afficient operation. The CCN representatives in the protocol design process felt that fundamental differences in approach arose from the "paper tape" orientation of most of the designers (the XDS 940, PDP+10, TENEX tradition) as opposed to the "unit record" orientation of the IBM systems. For example, this difference made development of a usable File Transfer Protocol very long and difficult, and the result was never fully satisfying to anyone, The "paper tape" maghines concretly used ASCII encoding and character-atematime echopies operation of terminals, compared to the EBCDIC code and linewatwamtime operation of the 360/91. Finally, most ARPANET designers and systems programmers were heavily oriented towards interactive service, while CCN primarily offered batch processing.

A result of the dominant orientation of the ARPANET community toward interactive computing was that CCN had difficulty providing user access to its batch Processing capability. A number of potential user groups might have accessed CCN through the ARPANET (f they could have dialed in with their batch job entry terminals. However, meither the common user Hosts (e.g., TIPs) nor the mejority of the other service. Hosts provided access for the synchronous types of batch terminals which are in common use. One solution to this problem was to make the IBM time-shaping system TSO available via TELNET, allowing remote users to create job Stream files at CCN and submit them for batch processing. A second solution, for a quer Host system, was to program a NETRJS user prodess (see Section 2.5.5). This happened slowly, but eventually seven different Host systems (mplemented NETRJS user programs, as shown in Table 3=1. However, since NETRJS was a "nonestandard" protocol, CCN never felt entirely comfortable in basing its user services upon the implementation of NETRJS et other Hosts.

With the exception of the official RJE protocol and graphics, CCN attempted to implement completely every user-level protocol. An outstanding example is CCN*s File Transfer Protocol (FTP) Server, which is the most complete implementation of that

RANO

OS/MVT on 370/158
(originally on 360/65)

UCLA=NMC

SEX on Sigma 7

Illinois

ANTS on PDP=11

MIT=DMCG

Harvard, Utah

DEC 10/30 on PDP=10

UCS8

Table 3-11 Sitem Supporting NETRUS Use (Seeden73)

TENEX on PDP=10

131, 89N, NIC, 14

protocol on the ARPANET, since CCN believed that efficient and reliable transfers of large files would be very important to its users. The effort was largely frustrated by the generally minimal efforts invested by most other Nosts. For example, no other site has implemented the restart gapability or the highwefficiency block modes.

It is obviously impossible to identify all the uses made of the CCN 360 via the ARPANET,. Some of the major computing efforts which were performed are:

- The Rand Corporation's meterological group, CCN's first large user, operated an atmospheric circulation model, which included programs requiring 850K and 1200K bytes of memory, In four CPU minutes, the 360/91 simulates six hours of real time. It was estimated that the same generation would have required 40 CPU minutes to complete on Rand's 360/65. Since their usual practice is to simulate 90 days of real time for their runs, it can readily be seen why they needed access to the bigger, fester machine.
- At the University of Illinois, both the Laboratory for Atmospheric Research and the Center for Advanced Computation make use of CCN's 360/91 through a variety of terminals and output devices at their computer center. These include a Gould electrostatic printer, an IMLAC programmable terminal, a Computek storage tube, and a penatype data plotter. As part of their work, they have produced displays which present a three-dimensional plot of pressure distribution within a cloud.
- The ILLIAC IV users group at NASA Ames. Research Laboratory has used CCN's 360/91 to test programs to be run on the ILLIAC IV. When the programs are successfully run at CCN,

they are then run on the ILLIAC IV, and the results are compared to validate the results produced by the ILLIAC IV.

- 4. Professor Hearn of the University of Utah used CCN over the ARPANET to do the formal algebraic manipulation of problems occurring in highwenergy physics. Professor Hearn ran his REDUCE processor to evaluate and expend algebraic expressions, using the 360/91 because of its speed and memory size.
- S. Teledyne=Geotech in Alexandria, Virginia, has used CCN for signal processing (data reduction) in support of seismic research. They were using their IBM 360/44 to perform fast fourier transform enalysis on their data. The data involved exceeded the capacity of the 360/44.

Overall; in spite of some feeling of alienation as a batch, "unit record", service center system in a time sharing "paper tape", research system world, CGNfs experience as an ARPANET server has been favorable. From a financial point of view, ARPANET usage of CGN computing services ranged from 10% to nearly of 20% of CGNfs income between 1970 and 1975, providing an important aid to balancing CGNfs budget. Obviously, CGN was still primarily a local service center for student and faculty

needs. However, from the Government viewpoint, CCN was an excellent resource. The 360/91 hardware was a half a generation beyond the 360 technology, and remained price-competitive with other dost recovery centers until recently. At ARPAFs suggestion CCN modified its accounting formulas to provide equitable charging for large-memory jobs; the Rand meteorological project proposed to run jobs requiring 2400K bytes for hours at a time; This change was useful to both ARPA and UCLA.

There have been many other benefits to UCLA from CCNFs involvement with the ARPANET. For example, programs installed at CCN by or for ARPANET users were also evaluable to all other CCN users, including EISPACK, SPEAKEASY, REDUCE, NASTRAN, and PL/M.

The ARPANET services that CCN performed were often related to, or at least similar to, existing faculty projects. For example, the Rand Climate Dynamics project was dependent upon the models and codes created by professors Mintz and Arakawa of UCLA; In another notable example of remote collaboration, Professors Donchin at the University of Illinois and Videl of the UCLA Computer Science Department used the ARPANET through CCN to collaborate on collecting and reducing experimental date.

The ARPANET provided a mechanism for several CCN staff members to interact with a community of bright system designers and builders, and provided the stimulus of new ideas and varied backgrounds. Furthermore, CCN menagement has observed that the ARPANET software design for the 360/91 was an interesting and challenging project, undoubtedly a factor in keeping several of the best systems programmers on the staff, Many programmers on the systems and user services staffs of CCN gained valuable exposure to the non-IBM world through ARPANET involvement, and gave the staff new insights into how much better software can be; This contact has hed many subtle benefits and reiged the level of system design and user interfacing at CCN.

3.1.2 TENEX Systems

The TENEX system is a virtuel memory system built eround the DEC PDP=10. The TENEX operating system, together with the associative memory page-mapping hardware which is used to support it, were developed at Boit Beranek and Newman, TENEX implements a virtual processor with a 256K word (35mbit words) memory for each user process. In addition, it provides a multiprocessor job structure with software interrupt gapabilities, an interactive command language, and advanced file handling capabilities. Files

are generally interpreted as character stream devices; thus real terminals. Network Virtual Terminals, and other network connections gan easily be viewed by user processes as files.

The TENEX system has proven quite popular within the ARPA research community, perhaps due to its early (about 1970) implementation of virtual memory at quite low cost, and its highly human-engineered user interface which grew out of the tradition of the XDS 940 Project Genie Software, This popularity is reflected in the fact that 17 TENEX systems are listed as servers in the 1976 ARPANET Directory, Further, nine of these systems are located in only two centers, at the University of Southern California's Information Sciences Institute (ISI) and at BBN's Research Computer Center (RCC). In fect, ISI and the RCC are the two major suppliers of TENEX service to the ARPANET.

Although TENEX is a general spurpose timescharing system, with the ability to run elmost all software written for the PDP=10, it is probably true that the vast majority of ARPANET use of TENEX is related to document production (editing, formatting), computer meil (reading, sending, filing), or artificial intelligence work using Interlisp (Interlisp is a dielect of the list-processing language LISP, the development of which has been

distributed among people at several ARPANET s(tem). Generally speaking. TENEX systems have been specialized to be either more suitable for text processing and mail or to be more suitable for tunning LISP.

Both the ISI and RCC centers are heavily involved in system software development, computer acience research, and networking experiments; thus each site initially appeared to place less emphasis on stability, high availability, and reliability than more service-oriented sites (e.g. CCN). Both sites approached the problem of integrating their research with provision of stable service by the adquisition of additional hardware to construct a "service" system and a "development" system, with the development system also available for beckup use. Once this step was taken, these two sites were able to offer significantly better performance than other, single system, TENEX sites and this fact probably contributed significantly to the gravitation of new users to these sites, supporting additional systems, achieving increased economy of scale in operation, and thus adding to their competitive mergin.

It is interesting that meither ISI nor BBN provides much user support in the form of consultants or hardcopy

documentation. The TENEX sites tend to supply small smounts of demputational service to very large numbers of users, in contrast to the service pattern typical of a batch center. For this reason, user support of the type offered by, for example, CCN appears to be prohibitively expensive. For this reason, the TENEX service sites have tended much more to online documentation and "consulting" via mailbox services.

In spite of the fect that the RCC and ISI are in some ways competitors for TENEX business, it is worth noting that they are quite different types of organizations and thus true competition does not really exist. The RCC is part of a commercial enterprise which has furnished its own computers and is attempting to make a profit in selling services. ISI is part of a non-profit organization which is operating government-furnished computers and is being reimbursed by the government for operational expenses. Therefore, direct government TENEX use tends to take place at ISI; other ARPANET use, for example by private contractors working on government projects, is somewhat more competitives.

3.1.3 Multion

Multica is a general-purpose computer colleborative effort by the MIT Laboratory (formerly project MAC), Honeywell Informative General Electric computer of vision) Laboratories. Multics was designed with being a prototype "computer utility", concepts and philosophy of serier times everal directions. The ability to share do the data sharing, is handled in a Multics structure.

Multics at MIT is a virtual memory sys Holos system with two CPUs, 384K word migrosecond core, en. 2 million words microsecond core). The hardware archit segmentation and paging, A hierarchy protection "rings" are used to limit relatively flexible way and provides a best build a secure (in the military sense Machine resources and accounting are u software that gives users who menege control over consumption patterns within th TENEX, there has been extensive human=engitarriace (although the interfaces are quit

During the stage of development of Host protocols, the Multice representatives, like those from CCN, tended to feel that they were seriously "outnumbered" by the representatives of the PDP=10 sites, and the TIP designers. The major points of friction included:

- 1) Multice was a linewate-setime system and poid a fairly stiff panalty in system overhead for fielding terminal input on a character-steetime basis. Nevertheless, the TELNET protocol design meetings were strongly oriented to character-steetime operation.
- Multics distinguishes between upper case and lower case characters; it is important to use the correct case in, for example, specifying a user name in a message. Most other systems, notably the PDP=10 systems, mapped lower case into upper case during the specification by the user of file or user names, as an aid to users with Teletype Model 33 terminals. Considerable pressure was put on Multics (not TENEX) by the ARPA office to "make it easy for Multics to accept TENEX=generated network me()."
- 3) Protogol for file transfers, especially for the transfer of mail, took a rather desual approach to access control (1,0),

authentication and system security). Because of the Multies emphasis on being a prototype computer utility, and the type of fairly stringent access controls which that implies, it was difficult for Multies to accept this casual approach.

Multime development and operation is marked by an organizational division which is rather unique in the ARPANET community, and which has led to less use of Multics via the network than would have been predicted based on the interest in, and capabilities of, the systems. The development of the ARPANET herdware and softwere was carried out by the Laboratory for Computer Science (LCS), with a staff of graduate students and researchers. The Multida Berver, however, is operated by MIT Information Processing Services (IPS). a service organization which has had difficulty in generating internal enthusiasm for use of this herdware and software in a production environment. At pertly for this reason. Multics has been less leest aggressively promoted then other systems. Nevertheless, (t is estimated that 15 to 20 persent of Multips usage is via the ARPANET.

It is worth noting that IPS has vigorously pursued the goal of insuring that local users do not indirectly subsidize network

users. This has led to a policy of handling as little of the protocol as possible at system level (i.e., interrupt level); this policy, in turn, has in some cases interfered with the capabilities which were intended by the protocols. As an exemple, performing HosteHost protocol flow control in the user's space, yet reading and writing the IMP interface in system space introduces an extra level of data transfer, and to some extent masks the effectiveness of the flow control mechanisms. Other network servers have been somewhat more released about distributing the overhead associated with the network among all users. Just as the cost of other I/O devices are distributed.

3.1.4 ILLIAC IV

The ILLIAC IV is a very large scale parallel processor, with a control unit and 64 processing elements. Each processing element has a local memory of 2K words (64=bft words); the control unit can address all 128K words. A considerable amount of secondary disk memory is provided in close association with the machine, and a very large UNICON laser storage device is included in the ILLIAC complex, elthough it is not yet operational. Actually the 64 processing elements represent one "quadrant" of the original design, but the remaining three quadrants were never bu(it.

Access to the ILLIAC is via a number of PDP=10 and PDP=11 computers which serve the role of peripheral I/O processors and front-ends. These, in turn, are sonnected to a number of dommunications systems, including the ARPANET. In fact, eithough the beginning of the development of the ILLIAC predated the ARPANET by several years, the concept of user access to ILLIAC via the ARPANET has been central to ILLIAC site planning since 1970; at times it was anticipated that all access to the ILLIAC would be via the ARPANET, and several different IMPs were installed to make this possible.

An IMP was first installed at the University of Illinoise Center for Advanced Computation to provide for ILLIAC access after the maghine had been delivered to that location. When the Illinois campus began to seem unsuitable as an eventual site for the ILLIAC, the seerch for a new site was conducted with the knowledge that users of the system would be able to access it vie the ARPANET, and thus the geographic locations of the system users did not need to have a first order effect on site selection. In early 1971 an IMP was installed at the Burroughst Paoli, Pennsylvania facility where the ILLIAC was being constructed; this IMP provided access to the 86500 which was originally intended as the ILLIAC I/O interface. Both an IMP and

a TIP were eventually installed at the Institute for Advanced Computation at the NASA Ames Research Center when this site was selected for ILLIAC installation.

ILLIAC connection, and the enticipation that all user eccess to both the processor and the associated "trillion laser Store would be wis the ARPANET, had some effects on the ARPANET which are not completely obvious. It was expected that huge volumes of data, requiring relatively high transfer rates, would be funnelled into the viginity of the ILLIAC: the direct result of this expectation was the initiation of the Plusibus IMP project, with a design goal of hendling connections operating at over a million bits per second. Two of the ARPANET inter-IMP sircuits terminating at the ILLIAC's IMP were upgraded from 50 to Kbs, and one of them (from the Moffett Field IMP site) was ordered to be further upgraded to over one megabit speed. Eventually, however, it was realized that much of the bulk transfer traffic to ILLIAC would be classified, and modes of access not involving the ARPANET were shown as a result.

From a system software standpoint, the ILLIAC is rather primitive; a user takes control of the entire system, runs his own software, and then relinquishes the system in a

single-programmed, batch-type mode of operation. A veriety of special languages including CFD, Glypnir, and Tranquil have been written to assist programmers in taking advantage of the computational parallelism provided by the system.

Use of ILLIAC and the associated facilities has been administratively limited to ARPA-authorized and NASA-authorized user groups. The customer support problems feded by the ILLIAC staff are thus limited to interestions with a relatively small and knowledgeable community.

It has been reported [Feik76] that NASA Ames researchers use about 20 percent of the ILLIAC's operating time (which totals 60 hours a week) solving two-dimensional serodynamic flow partial differential equations, and that they would like to process three-dimensional flows but have run into processor speed and memory size limitations, Other uses of the ILLIAC include global climate simulations from Rand and other NASA weather simulations, signal processing, seismic research, and linear programming; in general these uses are carried out via the ARPANET.

3.1.5 The Datecomputer

Like the ILLIAC and the 368/91, the Datecomputer serves as a specialized resource to a widespread community. However, rather then specialization as a computing engine, the Datacomputer is specialized to be a large-scale data management and storage utility.

more than the servers discussed previously, the Datacomputer was designed specifically for use in the ARPANET The Detacomputer designers at the Computer env(ronment. Corporation of America (CCA) acted on the belief [Maril175] that within a resource sharing network there is a natural tendency toward specialization of network nodes, such that the factoring of problems into their constituents, the assignment of these constituents to the appropriate machines, and the recombination of results will tend to become an automatic process. In the iimit, specialized network nodes become what may be "utilities", that is machines which perform a restricted range of functions solely for the benefit of the other machines. The Detadomputer is a network utility in this sense. It is entirely specialized for the performance of data management and storage functions. The Datacomputer designers further epeculate that the trend toward specialized network utilities will continue, and that the traditional standwalone general purpose machine will

evantually disappear from the scene. The computer world envisioned in such a speculation might consist of a network containing a few very large Detecomputer=like systems, a few very large computational utilities ("number drunchers"), and a large number of small human=interaction units (such as intelligent terminals), having limited computational power and local storage.)

Physically the Datecomputer consists of a processor, three levels of data store, and a connection to the ARPANET. The processor is a PDP=10 TENEX system, with primary core memory storage of 336K words (36=bits), Secondary atorage consists of eix apindies of IBM=2314=equivalent disk (24 million words total) and four spindles of IBM=3300=equivelent disk (68 million words total) which are used as a staging area for tertiary storage. The tertiary atorage device is an Ampex TerasBit Memory (TBM) system; which consists of control units built from DEC PDP=11 a to vilend pribactor a title series east cell a coccessor medebit per square inch. In the CCA configuration there are four tape drives, each of which contains 50 billion bits of storage, and transfers data at 5 magabits per second. A high speed seek from one and of a tape to the other takes approximately 45 seconds. The TBM can be expended to a total of 64 tape drives for 3.2 trillian bits of aneline storage.

The primary anticipated advantage of specialized network utilities is economy of scale, and the Detecomputer project fits this model well. The base cost of a TBM unit is very large; the cost of the TENEX frontmend is not small. Yet the cost per bit of the TBM is about \$1 per megabit, about a factor of 20 down from the cost of disk storage. Thus, while few individual installations could afford the price of a TBM, the pooling of many users requirements allows significent savings for all. In addition, of course, the cost of software development, maintenance and improvement can be spread over a large number of users.

The existence of the ARPANET motivated the development of the Datacomputer; the characteristics of the network shaped that development. The heterogeneous computer population of the ARPANET lead to extensive provision in the Datacomputer for a variety of byte sizes, character sets, numeric representations, and data stream structures. In most cases, automatic synthesis of, or conversion between, the elternatives is available. For example, one user could input a stream composed of variable length EBCDIC character strings terminated by a delimitor character. Another gould extract the same data as a stream of variable length ASCII strings each preceded by a length count,

The Datacomputer communicates with progrems that run on remote machines. The fact of remoteness prediudes the use of simple subroutine dells or similar meens of communication conventionally used within a single machine. The communication, furthermore, is not with people at terminals, who can be expected to make intelligent responses when failures unusua! 0.7 occur. but with programs. dircumstances Hence, all synchronization messages, error messages, language statements, and file descriptions must be creatable and readable by programs; likewise, a facility for checkpointing by user programs is required.

Detailanguage is the language in which all requests to the Datacomputer are stated. Detailanguage includes facilities for data description, for database greation and maintenance, for selective retrieval of data, and for access to a variety of auxiliary facilities and services. A besig characteristic of datalanguage is that all data is described. Descriptions are stored in the Datacomputer directory and are available to the user program in machine-readable format. A description contains the information needed to interpret the data, that is, information on data representations and structure.

The bandwidth available over the ARPANET, generally less than 40,800 bits per second, is two orders of magnitude less than typical local disk bandwidth. This means that different strategies and facilities are called for in using the Detacomputer over a network than for a user accessing data stored on a local disk. In particular, it is necessary for the user program to be able to send self-contained requests and to synchronize with and keep track of the remote task doing the user's bidding at the Detacomputer.

As an example of a self-contained request, one might want to update an employee file to indicate an agross-the-board select increase of 10%. With most local disk systems one can read each record and write it back updated. This would be relatively slow in a network environment. However datalenguage is powerful enough to specify this sort of operation to be performed entirely within the Datacomputer, insuring no network delays for data transmission.

The user has extensive control over the internal format of the data for efficiency reasons and flexible access controls are provided. At each directory level down to the individual file, data access and control privileges can be granted or denied to classes of users. The classes can be specified in terms of the

Host number and the Host-Host protogo; socket number from which the user communicates with the Datacomputer. The more usual password and login name restricted user classes are also evailable, as well as combinations of these with Host and/or socket number.

Quite eside from the control of access by password and address, the Datecomputer environment provides a stronger type of access control to date then is usually possible in a general-purpose machine. By definition, a general-purpose environment allows the programs within it enormous latitude in the functions they can perform, and it appears that programs can often be written to direcumvent existing access regulation procedures by taking advantage of errors that erise in unexpected circumstances. Such hostile programs are sometimes able, without authority, to access date, delete date, or drash the system and prevent other users from legitimately accessing date.

In the environment of the Detecomputer, the situation is quite different, since the system is logically a glosed, dedicated, special-purpose box, which responds only to a limited set of commands and does not provide a general-purpose computing facility. A hostile user program cannot be run on the box because the box does not run user programs. The approach can

inherently provide atronger guarantees that programs without proper edcass authority will not be able to access or damage data contained in the Datacomputer,

The TBM device was not available for use until late in 1976, but major use of Datalanguage and the sesondary storage facilities of the Datacomputer began several years sarlier. Some of the most significant Datacomputer uses includes

- Seismid Date: The seismid network project is covered in more detail in Section 3.3. A significant aspect of the project is the transmission of real-time raw seismid data to the Datedomputer at rates of 7×12 kilobaud, around the clock. By the end of 1976 nearly 70 billion bits of seismid data (raw readings, event summeries, instrument status reports, and various historical files) had been stored. If stored on conventional disk storage, this volume of data would have required more than 85 spindles of 100 megabyte drives. As the seismid database grew, researchers began retrievels against it, initially to develop and test procedures for use in ongoing seismid studies.
- 2, ARPANET Subnetwork Statistics: The ARPANET Network Control Center has been an established user of the Datacomputer, storing statistics on performance and usage of the IMPs and

circuits which implement subnetwork. Usage had grown to 600 megabits by the end of 1976, with steady retrieval activity against the data.

In Host Surveys: Another well-established use of the Detacomputer is the SURVEY application, derried on in conjunction with MIT's Laboratory for Computer Science.

Current survey data on the status of Hosts on the ARPANET continued to be stored at the rete of about 10,000 probes per day.

The SURVEY database has been used not only by those interested in the data per se, but also by researchers in the Very Large Data Base project at MIT. They used the Datecomputer's processing of requests against 2 quarters' worth of SURVEY data to test their work in estimating the cost of quary processing in very large databases.

4. ERDAL Several of ERDA's national laboratories have begun investigations of the Datacomputer; most active were groups at the Argonne and Lawrence Berkeley Laboratories. The major project is installation of a diffratological database by personnel at Argonne. This database contains 16 files, each with a year's worth of hourly readings for some U.S. city; the data are used by a number of sets of programs which model energy usage in buildings and sommunities.

3.2 User Systems

Just as there are many kinds of server systems connected to the network, so there are many kinds of systems which provide access to the ARPANET for terminal users. However, the vast majority of such systems are either Terminal IMPs or are DEC POP=11°s running the ANTS or ELF terminal support systems. Each of these is described briefly below.

3.2.1 The Terminal IMP

Initial terminal aggess to remote services on the ARPANET was via connection to a local Host. However, it quickly became obvious that use of a general-spurpose time-sharing system as a simple access port into the network was quite expensive. In addition, there was considerable interest in accessing the network from locations where there was no local Host in existence, These two observations led to the development of the Terminal IMP (or TIP) [Ornstein72] to provide direct network access for certain classes of terminals.

Logically a TIP consists of two independent machines, a normal IMP and a terminal-handling mini-Host. Physically, however, these two logical machines are combined in a single

Honeywell 316 computer identical to the 316=based IMP, but with the addition of 16K words (16=bit words) for terminal=handling program and buffers and a sophisticated Multi-line Controller (MLC) which provides the common logic for 63 bit-serial, character asynchronous terminals.

The TIP was designed to provide an inexpensive, reliable method for direct network access, and as such attempted to meet rather limited goals. First, it was felt that the TIP should serve human-operated interactive terminals, implying little need for flow control or error checking mechanisms between the terminal and the TIP. Second, it was assumed that alaborate terminals, such as graphics displays, would either contain a Host or access the network through a Host. Third, any computational requirement that might arise should not be handled in the TIP, but rather either by "intelligent" terminals or by the remote Hosts. The resulting TIP was therefore intended to provide access for users with "simple" terminals.

compared to normal computer installations. TIP (natallations are a bit peculiar. The TIPs are delivered without any systems documentation for the software system; in fact, since the TIP is part of an IMP, site personnel are forbidden to touch the system.

lest they leoperdize the integrity of the subnetwork. As with IMPs, problems are headled by the operators at the Network Control Center, Furthermore, since there is nobody locally who can be asked to modify the TIP system eccording to local desires, all requests for modification must be submitted to the TIP development group at Bolt Bersnek and Newman.

One type of continual pressure for modification is the steady introduction of new terminals to the market, and the desire on the part of the TIP user community to attach these terminals. In general, such terminals are "Teletype compatible", but each tends to have slightly different characteristics, for example in timing of print head movements to a new line, beginning of line, etc. Although differences like these may seem minor, each needs its own little bit of code to work properly or even to work at all. The general approach in the TIP is to provide each type of device with a piece of straightmline program to handle it. This maximizes TIP throughput by requiring the TIP to perform only the necessary perscharacter processing for any given device. The TIP actually handles Teletype Model 33%s, ISM 2741s, Odes line printers, Memorex line printers, Execuports, and any terminals compatible with them.

At one point after the TIP had been implemented, it was thought desirable to have a magnetic tape drive option for the TIP. Since this was clearly a device which could not be handled by the MLC, a standard Honeywell drive was used, and an ad hoc protocol was implemented to transfer tapes between TIPs so aquipped. However, it was impossible to smoothly integrate the tape equipment into the TIP structure, because its characteristics differed as widely from those assumed during the original design.

Other services for which there has been demand are things such as elaborate device status information, onwline server schedules, "news" and "complaint" facilities, and even access control and accounting facilities. It is clearly beyond the capacity of a TIP to handle these sorts of features by itself, However, since these are the capabilities normally associated with a service Host, it was logical to seek to provide such facilities at a larger Host and construct the mechanisms which dould allow a TIP user to gain access to them in a relatively transparent menner. The actual implementation of this concept makes use of the RSEXEC (see Section 3.5) on several TENEX systems (the reliability of the service is thus increased through redundancy), and an ad hoc "broadcest Initial Connection

Protocolⁿ which asks for a connection from each instance of the service and accepts the connection from the most responsive [Cosei175]. The ability to choose the proper location to implement a new function, as exemplified by such facilities, is clearly one of the adventages of a network environment.

It now appears that the original assumption that computing power should be in the terminal and not in the TIP is an oversimplification. Although it is reasonable to claim that a terminal access computer should not perform rotation calculations for graphics terminals, a terminal with a computer (such as an Imlas) would prefer to follow a more apphistizated communication protogol than is sufficient for Teletypes, For instance, the terminal may not always be prepared to accept characters when the TIP sends them and, conversely, the TIP is frequently not able to accept characters at the rate the terminal can send them. The MLC, designed for interactive terminals, does not lend itself to sophisticated flow control. One cannot merely say that the computational power should be in the terminal and not in the TIP; the very fact that the terminal has extra computational power requires the TIP to have some degree of matching power. However, the addition of relatively little computational power to the TIP can permit a large amount of additional computational power in the terminal.

3.2.2 ANTS

ANTS is an agronym for "ARPA Network Terminal System", a system which was developed at the Center for Advanced Computation at the University of Illinois. The ANTS design is much more ambitious then that of the TIP; it was developed to provide network access support not only to simple character terminals but also to devices such as synchronous batch terminals, card readers, line printers, tape drives, and sophisticated graphics display terminals [Souknight73]. The separation of terminal support functions from the network node permits additional capability in the system in the form of a disk, a wider variety of peripherals, more memory, etc., and thereby provides a limited amount of on-site processing power for the user to do housekeeping functions such as natwork accounting, onesite card desketo-printer listings, data storage on magnetic tape and disk, etc. It also provides him with a higher level interface into the various protocol streams than is provided by the TIP.

In addition to the terminal support and local housekeeping functions, the ANTS design was intended to offer a significant alternative to direct interfecing of large somputer systems to the ARPANET by acting as an "intelligent interface". An

"intelligent interface" is one which can be programmed to interface between dissimilar systems and map the input-output characteristics of one to the other. In such an instance, ANTS might be used to map terminal data streams into a form expected by the data-communications front-end of a large computer system such as the CDC 6688.

Internally, 411 devices which are available for use are grouped in classes for control purposes. In a typical ANTS configuration, there are several more or less standard classes!

- 1. Network connections
- 2. Files
- 3. Terminals
- 4. Non-file Structured peripherals
- 5. "Intelligent interfece" ports

Each class (or in some cases sub-class) is controlled by a device managers in ANTS are "virtual machines". They receive and execute "instructions" for passing data to and from their associated devices and the rest of the system. Other "instructions" provide for assignment of devices to system or user tesks.

Internally, all data is moved through the system in strings of 8-bit bytes called messages. Messages are transmitted along simplex "data paths". Each data path has a capacity for a

specific number of messages and a specified number of bytes which may be in transit over the date path at any given time. Thus a data path enforces implicit flow control. A second important capability of the date path is performance of data transformation operations. Such functions as code conversion, record compression and expansion, or aggregation of many small messages into large composites can be dynamically inserted in the data path.

The first ANTS system went into prototype operation at the Center for Advanced Computation of the University of Illinois in late 1971. By late 1973 it included the following hardwares

PDP=11/20 CPU 28K words core 256K word head-per-track disk ROM bootstrap Real=time=clock Interval timer Paper Tape Reader/punch 4 DECtape drives 2 9-track magnetic tape drives 36" drum plotter GOULD 4800 Printer/Platter Card Reader 2 Computek Storage Scopes IMLAC Grephics Display 8 Dialoin interfaces (110-300 Baud) 12 2400 Boud CRT terminals 2 TTY37 Teletypes

The ANTS project obviously was undertaken with a rather ambitious set of goals, siming to support local editing, the handling of complex terminals, and the "intelligent interfacing" of complex Host computers, all with a small Prodessor. Additionally, in order to make it possible for user organizations to modify and maintain their own systems it was decided to implement ANTS in an ALGOL=like language designed at University of Illinois apecifically for the purpose of implementing ANIS. Thus, the programmers and system designers were faced not only with the task of trying to understand and implement a large number of extremely domplex interfaces while maintaining high program bendwidth, but also the problem of using new language and its compiler. For reasons which probably included these, the ANTS project ren into serious difficulties, performance. Because of these both with achedules and difficulties, only a small number of the potential ANTS user sites actually installed the system on an operational basis.

3.2.3 ELF

ELF is a multiprogrammed operating system for the DEC PDP=11 (ELF is German for "eleven", and is recognizably similar to IMP) which currently has much the same goals as ANTS (Retz75)

[Retz76]. However, it was originally designed to provide a multi-user interface to computing resources via the ARPANET for solving problems requiring interactive signal processing depolities.

ELF development was started in early 1973 at the Speech Communications Research Laboratory in Santa Barbara, California, In early 1974 that system was tested at several network sites and a decision was made to expend the system depablishes. An initial version of the expended system was operational at the beginning of 1975. It is somewhat (ronic that the decision to expend the system to provide more of the (intended) ANTS capabilities was made at least partly because the ELF system became operational so much more quickly than ANTS, but a major reason for the large difference in the rate of development was the huge difference in initial goals.

The ELF system has a hierarchical structure. At the center of the system exists a set of modules, collectively referred to as the kernel, which perform tasks of resource management in a multiprogramming environment. An ELF process may be thought of as an autonomous stream of instructions (i.e., a program), having an associated program counter, general-purpose registers, stack

Processes may be in one of two states, ready or blocked. The kernel provides a set of primitive calls for outerelevel process, performing such tasks as the creation of processes, process synchronization, storage allocation, and sharing of an interval timer.

System procedures outside of the kernel perform tasks such as the interpretation of user gommands from terminals, maintenance of the system file structure, and control of communication with the ARPANET. The portion of the system which provides terminal support in ELF is referred to as the EXEC (or Executive), and is patterned after the user interface provided by the TENEX operating system. The ELF Network Control Program (NCP) controls the communication between processes in ELF and those processes residing in remote operating systems on the network.

The ELF system is being used in a veriety of applications which require a set of operating system components for network communication. The ELF Kernel provides a base for a variety of software configurations which are tailored according to system requirements. The Kernel, EXEC, and NCP modules support a number

of users who access the network by running the TELNET subesystem of the EXEC.

Additional processes may be added to the system for peripheral support. For example, a simple process which uses the ELF NCP to await a remote request for connection may be included in the system. The process accepts a connection, receives a stream of data, and outputs it to a local line printer. A status indication from the NCP signals that a remote process has closed the connection, and the peripheral control process then returns to the "listening" state.

The support of realetime date acquisition functions for speech research has been one of the goels in the design of the ELF Kernel. An EXEC subsequent performs functions of realetime data sampling and file I/O for digitization of speech waveforms. The system may thus be used as a data acquisition station while simultaneously providing a terminal support function for access to remote systems.

The ELF design has placed much more emphasis on development and maintenance via the ARPANET than ANTS. In addition, partly because of the limited initial ELF objectives, ELF is written in a macro assembly language rather than a highwievel language.

This admittedly makes modification and expension more difficult, but did probably make it easier to meet the resisting signal processing requirements. Crossepompilers for the BCPL and L1011 languages which generate PDP=11 code are evaluable on network TENEX systems for production of ELF user=level code.

Several debuggers which run on a stand-elone PDP-11 have been modified to run in an ELF waer address space. An additional debugging mechanism involves the interpretation of debug commands reprived from the network by a system debugging process. This minimizes space requirement in a processor being fissordqs. debugged while taking advantage of facilities available on systems to provide a comfortable user language and a symbolic representation of addresses and instructions. A system debug process which resides in ELF utilizes the breakpoint signal facility provided by the Kernel and is responsible for debugging a number of other processes. ELF system software is distributed as a set of source modules which are accessible in a directory at one of the network TENEX service sites. Users may access the source files by means of the network file transfer protocoi.

The ELF system has been quite successful in the ARPANET, both as a terminal support system and as a front-end. By late in 1976 there were between 20 and 30 ELF systems in operation in the network.

3,3 Specialized Communities of Interest

This section discusses three of the hundreds of specialized user communities which make up the patterns of network use. These are the climate dynamics research community, the seismic research community, and the "secure" user community; they are of special interest for quite different reasons. The climate dynamics research, carried out primarily under the direction of the Rand Corporation, is interesting because it used several large computers and was one of the first major activities making serious use of the ARPANET. The seismic research activities are of interest because of their heavy recletims component and their heavy demands for network bandwidths. Finally, the "secure" use of the ARPANET is of interest for the techniques used to make it possible.

3.3.1 Climate Dynamics

The Climete Dynamics Program at the Rand Corporation studied the effects on climate of large-scale changes in the environment by running global simulation models of the oceans and atmosphere. Experiments might study the effects of changing the mean temperature of an ocean, melting a polar ice cap, or placing a large lake in the middle of the Sahara Desert.

A typical experiment scenario was to first run a simulation of 90 days of the unmodified environment as a control standard, followed by a 90 day simulation beginning from the same initial data with the exception of the desired environmental modification. Two of the programs which are part of the modeling system require large regions of core, about a million dabit bytes, and have huge computational requirements. The model required about 40 minutes on the Rand 360/65 to simulate six hours of real time; the same simulation requires only four minutes on the UCLA CCN 360/91. Thus, the first requirement of the program was to use the ARPANET to access the 360/91 and later, when it become evaluable, the ILLIAC IV,

However, vital as access to the large processors was for running the simulations, the greater part of the programming effort was involved in the posteenalysis and display of the large data files that were generated as results. One 3deday simulation provides a full IBM 2314 disk pack, 28 million bytes of data. To be able to analyze this data meaningfully, a meteorologist must be able to soon the resulting weather picture and focus in on great of interest. He needs to look at such diverse items as the shange in temperature or pressure at a given point over time; or at the effect of an air pocket over the Pacific on the pressure

at the east coast of the United States; or on the humidity in Seattle. The postmensives were primarily performed by large betch lobe, but elso involved some interactive jobs. displays were large printouts (198 or so pages) and static graphics (magnetic tapes were generated at UCLA and sent to outside vendors for the generation of m(prof()m and prints; capabilities not available on the ARPANETI - 8 5 w=11 interective, over-the-net graphics from both ILLIAC and UCLA to a TEXTRONIX terminal at Rand. The display techniques used were not dynamically interactive, but did allow the researcher to choose the date and describe the display to be presented. The running of the graphics programs was a result of the connection of the terminal to the network via a 378/158 at Rand which not set up operationally to handle interactive user applications, The Rand computer was used primarily for storage of source programs, text editing, and listing output. It was felt editing functions were best kept logal to avoid interference with these functions when the network or remote servers were down.

Batch Jobs were submitted to UCLA via the NETRJS protocol; some result retrievel from UCLA and all interactions with the ILLIAC system used the File Transfer Protocol, However, since no site except UCLA (including Rand) (mplemented the compressed mode

of data transmission or the restart markers specified in FTP, there were annoying problems because of the amount of time involved and unreliability of the Hosts [Mobley77]. When a transmission "broke in the middle" the only recovery was to restart at the beginning. Another difficulty for the weather researchers was the non-standardization of dommon commands (LOGON vs SIGNON or LOGOFF vs LOGOUT) from one system to another, and interpretation of responses, especially when there were problems with the remote Host or its network interface. In addition, the network graphics protocols were insufficient for this group, which defined a standard low-level softwere interface for each graphics device, so that existing higher level routines could remain the same.

Nevertheless, the researchers involved in the climate dynamics project report that "we made extensive use of the ARPANET and the advantages were tremendous, It was an RJE facility. It was a connection to remote interactive systems. It was our means of moving data between remote computers. It reduced turn-around time and increased the efficiency of our use of programmers, time a great deal [Mobjey77]".

3,3,2 The Selamic Community

The seismic data activity involves the collection, storage and processing of seismic waveform information (seismograms) as measured by seismometers installed throughout the world. The data will assist seismologists in exploring techniques for detecting seismic events, pinpointing their location, and recognizing the causes of these events. A major application of the work is the detection of underground nuclear tests in preparation for future Strategic Arms Limitation Treaties. By establishing an on-line, real-wtime database of seismic information from a world-wide network of monitoring sites, a great deal of data can be made easily available to computers in the network for seismic analysis and other purposes [Dorin77].

The computers at various sites involved in the gathering and subsequent analysis of seigmin data are known as the Vela Seigmological Network or Velanet. Many of the computers are on the ARPANET, which therefore was phosen as the most appropriate communications medium available for much of the system, The Velanet consists of two sites sending seigmic waveform information in real-time, the Large Aperture Seigmic Array (LASA) in Montans, and the Norwegian Seigmic Array (NORSAR), LASA data

is transmitted via leased telephone lines to an intermediate processor (the Communications and Control Processor, or CCP) at the Seismic Data Analysis Center (SDAC) in Alexandria, Virginia, NORSAR data arrives at SDAC via a satellite communications link of ARPANET.

Deta which is not transmitted in real-time arrives at SDAC on magnetic tapes from the Iranian Long Period Array (ILPA), as well as from other instrument clusters throughout the world. Non-array seismic data is sent by magnetic tape from various locations around the world to the Albuquerque Seismological Laboratory (ASL) of the U.S. Geological Survey. Both the real-time and non-real-time data arriving at SDAC, as well as the data concentrated at ASL, are forwarded through the ARPANET to the Datacomputer (see Section 3.1) at CCA. The seismic data traffic to the Datacomputer is 12 kbs around the clock or approximately 30 billion bits per month. There are plans for additional sites which may boost the traffic volume to 35 kbs.

Progessors throughout the Velanat can retrieve the seismic data. Progessors at SDAC, Lincoln Labs Applied Seismology Group (LL=ASG) and possibly elsewhere will be used by seismologists for this purpose.

There are two basic categories of files stored on the Datacomputer for the selemic project. First, the rew data and status files provide complete selemic readings at various instruments and arrays around the world, as well as associated information on the status of these finstruments so that the raw readings can be properly interpreted. Second, there are the derived event summary and associated selemic waveform files, These consist of a processed distillation of the first category into likely sejemic events, and the time segments of raw date that are associated with these events, The derivation of the second datagory of files from the first is performed by processors at the Sejemia Data Analysis Center.

The real-time data from LASA and NORSAR (a, as noted above, routed first to the CCP and from there to the Datacomputer, Processing and buffering constraints in the arrays and their associated computers, the ARPANET, and the CCP provide a very short interval of elasticity (less than 30 seconds) between the time the data is generated and the time it must be accepted at CCA. This small interval of elasticity has not only network implications, which are discussed below, but also implications for the operation at CCA.

The Datacomputer system is implemented within a general purpose time-sharing system with a verying load depending on time of day and other factors. It cannot guarantee the required responsiveness. In addition, both the basic computer system and its peripherals, including the TBM, require regular preventative maintenance, and hence cannot be operated continuously.

To provide the required roundetheeclock responsiveness, a small, dedicated, reliable system known as the Seismin Input Processor (SIP) has been implemented. The SIP is a PDP=11/40 with two ISM 3300=like disks. It accepts the data stream from the CCP, buffers it on its disk, reformets the data, and periodically transmits it to the Datacomputer. At the current bandwidth, 26 hours of buffering are provided per disk pack. Thus the SIP completely isolates the resisting data stream from Datacomputer downtime or dalay.

The SIP uses standard Host-Host protocol to communicate with the Datacomputer but uses a special protocol for the real-time path to the CCP. This special protocol eliminates the normal Host-Host handshaking and connection setup overheads and results in simpler, more efficient communication. Furthermore, the special protocol eliminates the standard protocol requirement

that no more than one message be in the network in one direction at a time. This modification increases bendwidth and decreases network blocking. The special protocol also maximizes network efficiency by packing logical messages into full size physical messages and uses sequential message ID numbers.

Seismic networking activities have placed considerably more strain on the operational ARPANET than most other uses. For obvious economic reasons there has been only one trans-Atlantic path included in the network, and this path is rated at only 9.6 kbs rather than the usual ARPANET rate of 50 kbs. Even more significant (a the IMP-to-IMP acknowledgement strategy which normally allows only eight packats to be outstanding on a circuit. Since the trens-Atlantia circuit is provided by a satellite with a minimum roundetrip delay of about one-shalf second, this strategy restricts traffic on this circuit to about is peckets per second, regardless of the packets! sizes. These factors combined to make it difficult or impossible for the Velanet to maintain the necessary flow rate from NORSAR to SDAC. One factor revesied by the analysis of the difficulty was that, because of the 16 packet per ascend restriction, five TIP terminal users in Norway and England, each typing one character per second and operating in remote echo mode, would saturate the channel (although achieving only a 64 bit per second data rate), Thus this usage pattern, if it actually existed, would be likely to reduce the Velanet bandwidth to unacceptable levels. This analysis lead to a modification of the inter-IMP strategy to allow more outstanding packets over the satellite channel.

The mode of operation at CCA also led to intra=IMP bandwidth problems. The SIP needed to continuously absorb about 12 kbs of traffic from the CCP, and from time to time unload accumulated data to the Datacomputer. A ressonable data rate for this latter transfer is on the order of 100 kbs. At the same time, of course, other network storemandsforward traffic; and other Datacomputer users, must be accommodated by the CCA node, The TIP originally installed at CCA had insufficient internal bandwidth to meet theme requirements while also providing terminal aupport. The TIP was first replaced with an IMP, but the IMP was required to support a VOH donnaction which usurped about one-third of the normal IMP buffering, so in spite of marginally sufficient bandwidth the IMPfs internal and to end protocol buffering requirements resulted in continued Unsatisfactory performance, Eventually, the VDH connection was moved to another network node, and finally a higher-capacity Pluribus IMP was installed, (The ARPANET's first Pluribus IMP

was installed at SDAC to meet the requirement for more external connections, Hosts and inter-IMP circuits, than the Honeywell IMPs could handle as well as the anticipated bandwidth requirements.)

Yet another effect of the Velanet use of the ARPANET has been the necessity to pay special attention to the large=volume seismid data flows when designing the network topology. For example, because the requirements for source=IMP and destination=IMP buffering to support a given bandwidth increase in proportion to end=to=end delay, and delay ingreases in proportion to path length (number of hops), it is desirable to have path lengths from NORSAR to SDAC, SDAC to CCA, and ASL to CCA as short as possible. These considerations resulted in the relocation of a previously=existing Washington to Soston directito SDAC at the Weshington end and to CCA at the Boston end. A Gunter to Mitre circuit was also partially motivated by Velanet needs.

3.3.3 Secure ARPANET use

As the ARPANET moved from a computer communication research tool to an operational service device there began to be interest in using it to corry classified data in a secure manner. The

uncontrolled ("black") environment is "end to end encryption"; the data passes through an encryption device (a Key Generator) as it enters the communication system and through a decryption device as it leaves. This technique is not directly suitable for a packet network such as the ARPANET because each IMP needs to process a portion of the data (the destination address). A solution to this problem in military message switching is known as "link encryption"; the message is encrypted as it besses into each communication circuit and decrypted as it enters each switchs. This solution, however, is only acceptable if the switches themselves are in a secure ("red") environment, glearly not the case in the ARPANET.

The initial solution chosen is a device known as a Private Line Interface (PLI) which acts as an interface between a secure Host and the non-secure ARPANET. The PLI drives the Host-supplied data through an approved Key Generator (KG) unit. A fixed address, supplied by the operator at PLI initialization time, which is unencrypted ("clear") is then appended to the encrypted data and delivered to the IMP.

In detail, the PLI is a Pluribus computer which appears to the Host to be an IMP, and to the IMP appears as a Host. The PLI is really two computers in series (See Figure 3=1). The secure Host communicates with the redehalf PLI via a normal Host-IMP interface. The red-half PLI signals the KG unit to generate and send a key-sequence. The black-half PLI supplies clocking to the KG. It agans the incoming data stream for a key-sequence, compresses the key, and stores it as the first few bytes of an ARPANET message. Then the redubalf sends a message segment through the KG unit. The segment is pedded with SYN characters if needed to bring it to a fixed size. The black-half adds this data block to the compressed key, assigns a message identifier and a message header (including destination) then transmits it to the IMP. It also notifies the redwhalf of the message identifier assigned wis a uniwdirectional data and status link (from black to red). When a RFNM or other status indicator is received from the IMP the redwhelf is notified over this same data 1(nk.

When the black-half receives an encrypted message from the IMP, it starts the KG, expands the key-sequence in the message, then transmits the data through the KG unit to the red-half, Since the black-half provides a clocking signal, the red-half must be prepared to accept the decrypted date as fast as it is

Section 3

DRAFT Expertence

Figure 3-1: PLI Configuration

sent. There is no way for the redehalf to indicate overrun or other errors to the blackwhalf, although errors are reported to the Host. This requires the Host-Host protocol to effect retransmission should errors of this type occur.

The first pair of operational PLI's were installed in the ARPANET in early 1976 in Sunnyvels and San Diego, California. As a result of the experience gained from initial operation of this pair, approval for some edditional control signals (Rad Receiver Full, Black to Red Word Empty, Reset) from the radehalf to the blackwhalf was sought from the National Sedurity Agency and later received. Approval of such signals is necessary because they contain the potential for passing classified data into the black environment without encryption through applicant, melfunction, or malicious design.

The most serious difficulty with the PLI design for the ARPANET environment is that a Host using the PLI is constrained to communicate only with the one other Host whose address is set in the black-half PLI at system start-up. This deficiency was accepted in the initial design in order to shorten the development and approval cycle, but with the approval and installation of the first pair of units, work began on the design

of a "multi-address" PLI. The intent of this device is to incorporate a low-bandwidth unencrypted data path from red to black which would allow a few address bits (5) to be associated with each encrypted data block. The bleck PLI could use these bits as an index into a short table of prespectived data.

The most serious obstacle faced by this plan is the need for NSA to be able to certify the correct operation of the redwhalf PLI program, so as to insure that classified data is not transmitted through the unemprypted path. The language chosen for implementation of this program is Meta-BCPL; this choice was made to assist NSA in their verification procedure in order to increase the probability of achieving certification.

At the same time, the redehalf PLI was expanded to be able to handle up to four Hosts. Thus, assuming finel approval of the system by NSA (expected in 1977), a secure Host sonnected to the ARPANET through the Multi-address PLI will patentially be able to address 32 other PLIFS, each with 4 Host, or a total of 128 remote destinations.

3.4 Network Meil

One of the most remarkable outgomes of the existence of the ARPANET has been the explosive growth of what has been called "network mail", "message service", and so on. This growth is the more remarkable because it was essentially unplanned, unanticipated, and mostly unsupported in its early growth.

Early in the course of ARPANET development it became obvious that the researchers involved with, or connected to, the network would want to communicate with each other in various ways. As part of its initial work in facilitating the communication between Most-Host protocol developers, the Network Information Center made use of the SRI On-Line system (NLS) available to this group; NLS [Englebart73] provided a relatively sophisticated text management system which assisted individual researchers to exchange notes, dollaborate on document production, create new text structures (essentially files and records) from fragments of previous structures, and so on, In fact, NLS was, although not described in those terms at the time, a serious "automated office" system.

Unfortunately, there were some rather serious drawbacks to ARPANET use of NLS. First, the system was optimized for use from

a very elaborate work station, with high-bandwidth graphic output and some rather esoteric input devices. Most network users didn't have the appropriate I/O devices, and their simulation on ordinary keyboard/printers was quite awkward. Second, use of NLS required logging-on to the single computer on which it was implemented, a system which was not always evaluable and which, in any case, had the dapadity to support only a modest number of simulateous users.

At a much more primitive level than the facilities of NLS, most of the interactive time-sharing computers of the late 1960°s included software for "linking" terminals together, and for implementing "mailboxes" where short messages (usually to or from the computer operators) gould be deposited or retrieved. These tools were constructed to help remote users obtain assistance when problems developed. Such tools were in use at least as early as 1965 on CTSS at MIT'S Project MAC.

The frustration of extempting to use NLS from remote ARPANET sites, combined with the general availability of primitive "mailbox" facilities on many of the network's Hosts, led relatively quickly to the development of simple programs which essisted a terminal user in creating outgoing messages and

reading received messages. The origin of these programs is uncertain, but has been attributed to the 88N TENEX system [Henderson77].

Network message service was an (mmediate success. Message flow grew in volume to become the most visible traffic component on the network. Use of the service has had a substantial impact on the organizations involved, stimulating dramatic shifts of dependence away from the traditional media (postal service, telephone).

 with principal investigators, protocol designers who needed to exchange documents, and a variety of other interest groups. This situation is probably different from, say, the population of users of a commercial time-sharing system or network.

Given the key differences between the ARPANET community and widely distributed sets of individuals with access to electronic transmission systems, as described above, there are several reasons why natwork mail came to dominate more traditional communications made. These include:

Transmissions speeds A message can be delivered from any imput point to any mailbox location in only a few seconds. This is a result of the ARPANET communication system design.

Senderwregelver decoupling: Although transmission speed is almost "real time" in rate, there is no need to establish a real-refere channel between the sender and the receiver as there is in telephonic communication (as well as the most domain modes of TELEX and facsimile transmission), Requirements for establishment of a real-stime channel are often time consuming for a caller and disruptive to the called party. The decoupling is a result of the buffering in the mailtox hosts,

- Location-independent reception: Since access to a mailbox (a via the network; a mail recipient (a not isolated by mobility so long as network access is possible. This is a result of the widespread geographic coverage of the ARPANET and of the availability of truly portable terminal devices beginning in the early 1970°s.
- Location-independent transmission: Just as it is a major benefit for a mail recipient to evoid isolation during travel; so it is a major benefit for a sender to be relieved of the burden of tracking his correspondents. This is a result of the fixed location of a mailbox (as contrasted with the telephone, TWX, or TELEX).
- Hultiple copies: Much interpersonal communication is usefully shared by more than two people; for example, a small group of researchers working on a common project. This fact is perhaps best examplified by the growth in office copier usage since its introduction. A piece of network mail can be sent to several meliboxes as easily as to one. This is perhaps the most important factor in the use of mail; unlike the telephone, it is easy to communicate among more than two people simultaneously, while unlike the postal service the

interactions can be at electronic speeds. This is a result of the oneline (electronic) "storage" of the message at its point of generation, which makes it trivial to produce additional copies.

The conscious (or perhaps unconscious) recognition of these factors was instrumental in convincing large groups of network users to begin communicating through the mailbox fadilities available on most of the time-sharing Hosts. However, as network mail began to supersede the use of postal and talaphone services It became more and more obvious that one wished to deal with network meil in the same ways that have developed over the course of time for traditional office documents. That is, one wanted the ability to have lists of primary recipients, secondary (carbon copy) recipients, "blind" copies, file copies, and so on; dating and authentication (a signature surrogate) were seen to be desirable; it was useful to be able to "write notes in the margins", to forward messages, to keep specialized mailing lists, to engage (n bulk mailings and so on. It also was seen to be useful to send messages to, and rem regetye messages from, programs; for example, there are many "desk palender" programs which send daily reminders of scheduled events. In other words, the widespread use of simple network mail ied inexprably to the

desire for more autometed office functions, as forecast by the developers of NLS; Thus, within a few years it seemed valuable to focus research on ways in which the network mail facilities could be extended to create a structured, multissite "automated office".

Initially measages were thought of as simple text objects, As message manipulation programs became available, structure was imposed upon these objects to separate out a collection of header fields. The obvious next step was to view the message as a whole a atructured object. Conventions for interpreting the structure permits senders to generate, and recipients distinguish, those components of the message which correspond to such classic message parts as the sender field. Thus the growth of more sophistidated mail processing programs led to the obvious necessity for mail formatting protocols (the mail itself is transferred from Host to Host via the file Transfer Protocoll, such protocols have been developed (on a relatively informal basis). Of course, the protocols not only specify the location, or method of identifying, the standard fields, but also the way to interpret the contents. For example, the "To" field should be thought of as a collection of addresses, the "Subject" field as a text string, the "Date" field as a detectime group, and so on,

More recently, some ARPANET message services have permitted extensions to the basic conventions by allowing message authors to define additional fields for their own purposes. Thus, communication on a fixed subject (for example, a company's contracts) could be facilitated by special message fields germane to that subject. Contract related messages might contain fields for contracting agency, start date, and date, and costs.

To enable a message program to serry out such operations as sorts and searches on these user fields, the detaitype of each such field must be made known to the program. Consequently, the agreed-upon structure of messages must be righ enough not only to support the existence of user fields but also to ellow the inclusion in each such field of the data type of its contents. Thus the intention of the author can be conveyed to the resipient and to the program through which he manipulates his messages. Simply stated, messages must be somewhat self-describing, and to be able to interpret the self-description on a variety of Hosts a standard protocol will be required.

There are, of course, many open (saues regarding network mail, First, although many investigators have asserted that some form of electronic mail will dominate the future of

communications there are a number of stumbling blocks. Of these the least are economic and technical; far more important are the opposition of established communication groups and regulatory agencies [Panko77]. However, additional technical development is still dlearly necessary to assure privacy and authentication, handle long documents (several printed pages) as conveniently as notes and letters, incorporate graphics (at least sketches and line drawings) in a simple way, and reduce the cost (dellars, weight, bulk) of convenient transmission/reception units (terminals) before it is reasonable to expect use of such systems by the average person for either business or personal communication.

3,5 Distributed Systems

The implementation of the ARPANET communication system created an environment in which it was possible to work on the development of multi-computer distributed systems. This section describes two such systems; the Resource Sharing Executive which was primarily implemented on TENEX systems, and the National Software Works which is a heterogeneous system. Also mentioned is the TIP Service Facility which uses the facilities of RSEXEC to provide TIP services. All of these systems are "network-based" systems, which not only use the ARPANET for communication among components, but wouldn't even have a reason for existence in the absence of a resource-sharing communication system.

3.5.1 The Resource Sharing Executive

The resource sharing executive (RSEXEC) is a distributed, executive-like system that runs on TENEX Host computers (Thomas73), By sharing resources among themselves the Hosts can provide a level of service better than any one of them could provide individually. Within the environment provided by the RSEXEC a user need not concern himself directly with network details such as communication protocols nor even be aware that he

is dealing with a network. The RSEXEC system was developed as an experimental vehicle to explore a wide variety of network operating system issues ranging from the types of features that would be useful to network users, to system structures and mechanisms for implementing those features, to strategies for ensuring reliable, fail-soft performance.

The goal was for RSEXEC to provide acress to the combined resources of the TENEX Hosts much as the TENEX operating system and command language interpreter (called the EXEC) provide access to a single TENEX Host. The RSEXEC was designed to serve both as a command language interpreter for users and as a program execution monitor for user programs.

The RSEXEC system has three principal components; the RSEXEC program; the RSEXEC Server (RSSER) Program; and the RSEXEC/RSSER protocol (see Figure 3=2); There is an instance of the RSEXEC program for each active user of the system. An instance of the RSSER server program runs on each of the Hostis; Its task is to make the resources of its Host accessible to remote users on demand by accepting requests made upon it by remote RSEXEC programs. The RSEXEC/RSSER protocol is the set of conventions governing the interactions between the RSEXEC and

Figure 3=21 Components of RSEXEC

RSSER programs. It defines a set of commands and responses necessary to support the various RSEXEC system features

Although primarily a homogeneous system comprised of similar TENEX Hosts, dissimilar Hosts (Multime, IBM 360/TSO, PDP+10/ITS) have been partially integrated into the RSEXEC system [Forsdick77]. The method of integrating a new Host type into the system is to implement RSEXEC and RSSER programs for the new

Host, The functional requirements of the RSSER program are specified by the RSEXEC/RSSER protocol which was designed to be independent of Host type. The philosophy for the RSEXEC program is that for a given Host type the program should provide local users and programs access to the combined resources of the RSEXEC Hosts in the same style that the local operating system provides access to local resources. Thus, a user's interactions with a TENEX RSEXEC program would be similar to his interactions with a single Host TENEX system, while a user's interactions with a multice RSEXEC program would be similar to interactions with a single Host Multics system. Although the style of interactions could be expected to vary from Host type to Host type, the functionality supported by RSEXEC programs for different Host types would be similar.

The udefulness of multi-Host systems such as the RSEXEC is, to a large extent, determined by the ease with which a user can manipulate his files. Because the Host used one day may be different from the one used the next, it is necessary that a user be able to reference any given file from all Hosts.

The file handling facilities of the RSEXEC were designated to:

- 1. Make it possible to reference any file on any Host by implementing a file name space which apans across Host boundaries.
- 2, Make it convenient to reference frequently used files by supporting "short hand" file naming conventions, such as the ability to specify cartain files without site qualifications.

Uniform file access is achieved by providing a network-wide file naming syntax which extends each single Host file naming syntax for full file pathnames by adding a Host name field. Thus, any file in the system can be specified by its full RSEXEC pathname which includes its network location. This full pathname is interpreted in the same uniform way regardless of the Host of which it is generated. Convenient file access is schleved by silowing the use of partial pathnames to specify frequently referenced files. These partial pathnames are interpreted with respect to the user's own private working directory maintained by the RSEXEC program. A working directory, in general, may span Host boundaries in the sense that it may detalogue files from several Hosts.

The working directory contains an entry for each file in each of the component directories specified in the user's profile. At the start of each session the RSEXEC adds current detail to the user's directory by gathering information from the server programs at the Hosts specified in the user profile. Throughout the session the RSEXEC modifies the directory, adding and deleting entries, as necessary. The directory contains frequently accessed information (e.g., Host location, size, date of last access, etc.) about the user's files. It represents a source of information that can be accessed without incurring the overhead of going to the remote Host sach time it is needed. The user's "profile" mentioned above is explicitly controlled by the user's who can add or remove components to control the set of Hosts over which, for example, partial pathnames are interpreted.

The user can take advantage of the distributed nature of the file system to increase the "accessibility" of certain files he considers important by instructing the RSEXEC to maintain images of them at several different Hosts. With the exception of certain special purposes files, the RSEXEC treats files with the same pathname relative to a user's directory as images of the same multi-image file. For example, the user profile is implemented as a multi-image file with an image maintained at every component directory.

file operations initiated by executing programs are handled by a technique called "encapsulation". The RSEXEC program intercepts certain file system operating system calls made by programs before they reach the local operating mystem. It interprets the file operations in the context of the network file system as described above. Operations referencing local files are passed directly to the local operating system, operations are forwarded to the appropriate remote RSSER program. An important benefit of this engapsulation technique is that programs written for a single Host environment need not be rewritten to operate within the network environment. Thus, the value of existing software, such as text editors and language processors, is significantly enhanced since the software execution environment includes remote as well as local data files.

Another important class of RSEXEC features are concerned with system status and interactions among users. These features are supported by commands which allow users to obtain information such as the current loads of the constituent Host systems (e.g., the number of active users, load factors, etc.) or the particular users currently logged into the various systems. In addition, a user can request that his terminal be directly "linked" to the

terminal of another user (who may be logged into sither the logal system or a remote system) in order to engage in an oneline dislogue. At the user interface remote links are initiated and terminated in the same way as local links. However, remote linking is accomplished by dooperation between the RSEXEC program and a remote RSSER program, whereas local linking is directly supported by the local operating system.

Since the RSEXEC was originally created as an addition to existing autonomously administered TENEX systems, both its conception and implementation preserve the prerogatives of each of the Hosts. Each Host perticipating in the system totally maintains its own descriptive and state information about the resources that it supports. Thus, any RSEXEC Host can continue to operate autonomously, although with diminished scope. Such autonomous operation can, however, ultimately lead to conflicting interpretations of global file names. No attempt is made to maintain a long-term system-wide data base of network resources. Any replication of resources (e.g. multiple copies of a file) is evident only when an environment is dynamically configured for a user.

The total autonomy of the RSEXEC Hosts has the further administrative implication that a user must register individually with each Host whose resources he wishes to utilize. This involves establishing valid accounts, acquiring sufficient resource guarantees, setting up appropriate access controls, and so on for each constituent Host. Once this has been done, RSEXEC facilitates the use of the resources on these Hosts. Thus, elthough the RSEXEC eventually helps the user cope with the distributed environment, the user is still forced to first view it as the collection of independent, and perhaps even competing, entitles which predeted the ARPANET.

3.5.2 The National Software Works

The National Software Works (NSW) is an ARPANETebased system (Robinson76) designed to provide programmers access to a large range of software tools, e.g. text editors, compilers, assemblers, and debuggers, which can be used in the software development activity. From the standpoint of the programmer or program manager, the NSW environment consists of numerous software development tools, running on a variety of computer systems which are geographically and administratively distributed across the country, but accessible through a single

access=granting, resource allocating monitor with a single, uniform file system.

Active work on the NSW development began in July 1974 under the Joint sponsorship of the Air Force and ARPA. Initial operating capability using NSW components residing on two different computers (IBM 369/91 and PDP=19) was demonstrated in July 1976. The NSW project was conceived of as a possible way to attack the high cost and poor reliability of software within DOD by applying the results of ARPA research programs, especially the ARPANET.

It was noted that the types of tools mentioned above, which span the software development process from design through implementation and checkout, have been available on a variety of computers for many years, yet they are seldom applied together in an effective way to support software implementation projects [Forsdick77]. One reason is that existing tools have been implemented for different computer systems, and programmers typically do not have access to the range of machines that house them. Furthermore, even if they did, programmers would have to master a variety of different Host systems and command languages, and deal individually with basic interHost incompatibilities to

use the tools. As noted, the NSW mystem addresses both of these problems.

The major components of the NSW are the frontwend (FE) systems through which the users access the NSW, the access granting, resource controlling sentral component called the Works Manager (WM); Foreman (FM) modules that interfece tools on the Tool Bearing Hosts (TSHs) to the WM, and communication protocols (MSG) that specify the communication links between the various NSW components.

A scenario of how NSW can be used illustrates the nature of the environment for software development it provides. The NSW evstem includes as T8Hs PDP=10 Hosts operating under TENEX, Honeywell 6000 Hosts operating under Multics, and an I8M 360/91 operating under D8. Among the user sommunity is a group developing software for the AN/UYX=20 computer, a smell computer that does not itself support a wide range of software development aids, and is generally configured for production use rether than program development. The NSW will support the adjacement each cycle for this group in the following way. Interestive text aditors that, run on TENEX and Multics are available for program preparation and modification. The source programs will be

written either in assembly language or a higher level language for the AN/UYK=20. These source programs will be assembled or compiled by language processors that run on the 360/91 and loaded to form executable AN/UYK=20 object modules by loaders that can run either on the 360/91 or on TENEX.

Interactive debugging for the AN/UYK=20 programs will be done within NSW using a AN/UYK=20 debugger tool which runs on a TENEX TBH. This interactive debugger is, in fact, a multi-computer tool. Part of it runs on TENEX and part of it runs on a microprogrammable computer, the MLP=900, which is connected as a peripheral to one of the TENEX TBHs. The MLP=900 has been programmed to emulate the AN/UYK=20 and operates under the control of the TENEX interactive controller/debugger module. After the several edit=compil=edebug cycles required to produce a debugged AN/UYK=20 program, the user can "export" his debugged 4N/UYK=20 load modules from the NSW for final checkout and operation on an actual AN/UYK=20 machine.

To use any tool, a user simply specifies the tool by name (e.g., the TECO text editor, the AN/UYK=20 compiler, the AN/UYK=20 debugger) and the NSW system starts an instance of the tool and connects the user to it. The user need not know which

T8H supports the tool nor the command language of the particular IBH. He need only learn the NSW command language. To support tool execution the NSW implements a type of distributed file system. All NSW tools utilize this single NSW-wide file system for their file references. When a tool attempts to access a data file (e.g., a source program for a compiler, a load module for a debugger) the NSW system ensures that the access references the correct file independent of its actual location. If the referenced file (s stored on another Host system, the file is automatically transported to the referencing TBH so that the access can be completed.

In addition, NSW may perform file transformations to deel with basic interHost file system incompatibilities. Thus, the user and the tools he employs need not doneern themselves with the location of data files, the details of the file systems on the various constituent TBHs, the movement of files between Hosts, or the data translations that may be required to make a file dreated on one Host type useble on another Host type, Rather, they interect directly with a single, uniform NSW file system.

Each active user has a dedicated FE process which acts as his interface to the NSW system for the duration of his asssion, The FE acts principally as a command language interpreter making requests upon other components as necessary to satisfy user commands. The MM is the resource allocation and access control module for the NSW system. All attempts to access NSW resources, such as tools or files, must be authorized by the WM. To perform its task, the WM maintains system-wide data bases such as an NSW file system catelogue, tool descriptor information, user password and account information. and so on. Interactions between MM processes and other system components occur on a transaction oriented basis. That is, the system does not dedicate a single WM process to each aptive user for the duration of the user Rather, WM processes are dynamically allocated (and deallocated) as necessary to support a user session,

When a user requests the start of a tool, an FM process on the appropriate TBH is allocated to the user for the duration of the tool session. The FM process provides an execution environment for the tool and controls its operation. This execution environment differs somewhat from the standard environment provided to the tool by the local TBH operating system. For example, when a "file open" operation is initiated

by a tool, the operation must be processed in the context of the entire NSW rather than that of the local Host operating system, The FM process responds to such an attempt by dooperating with a WM process to complete the file reference. The WM process consults the NSW file datalogue to verify the existence of the file specified by the FM, and the authorization of user and tool to access the file. Next, the WM sats to ensure that the file can be physically accessed by the FM/tool. In general, this may require movement of the file to the FM Host and possible translation of file data to a form usable by the tool.

Communication between the verious NSW system processes well as the allocation and deallocation of those processes is the responsibility of a component selled MSG. Two sorts of addressing are supported by MSG. Generic addressing is the means by which a process initiates a trensaction with another previously unrelated process. Generically addressed messages are used when any process of a given type is acceptable. The second type of addressing is specific addressing. It is used when the sending process must communicate only with a particular process.

The FM component is designed to provide two different tool interfaces to NSW. An encapsulation interface is provided to

permit software packages that were developed to run under a particular TBH operating system to be installed as N8W tools. This interface is similar in concept to the encapsulation meshanism used in the RSEXEC system as described in Section 3.5.1. The second tool interface to NSW provided by the FM is a direct interface which permits tools to explicitly dell NSW functions. Generally the FM must interact with other N8W system processes to satisfy these tool calls. This interface is provided for tools developed specifically for operation within the NSW environment.

The NSW file system is built upon the file systems of the constituent Hosts. The file systems of the Hosts are used simply as a storage pool for NSW files. The WM maintains a file catalogue which contains an entry for each NSW file. A file entry includes the NSW name for the file, the location of the file in one (or more) of the constituent Host file systems, and various file attributes. NSW file operations are logically centralized in that the central WM file catelogue is (slmost) always used to resolve file references.

When a tool opens a file, a copy of the file is moved from the NSW file storage space into a tool workspace maintained by

the tool's FM. The file manipulated by the tool is truly a copy in that any modifications the tool makes to it will not be reflected in the actual NSW file maintained by the WM unless the FM/tool explicitly delivers the copy back into the NSW file system.

Because of its modular structure, the NSW system architecture is potentially resilient to individual Host failures. Indeed the interecomponent protocols have been designed to make continued system operation possible in the presence of FE and FM/tool Host failures. In addition, mechanisms have been developed to recover files trapped in a tool workspace due to a IBH crash. When a crashed IBH is restarted, tool workspace: (n use when the crash occurred are preserved in 8 way that allows a user to later retrieve selected files. At present, the centralization of the WM is the principal weakness of the system from a reliability point of view. Because the WM and its central date bases reside on a single Host, the NSW system is vulnerable to failure of this Host, A multi-Host implementation of the WM function would permit continued NSW operation in the presence of WM Host failures. In addition, it would also allow the system to pracefully expand as the number of users grows by making it possible to distribute the WM load among several Hosts. The principal technical problem here results from the fact that the WM data base is logically centralized. To distribute the WM, the data base must be distributed. This requires the development of synchronization mechanisms to insure that the distributed parts of the data base are maintained in a consistent manner.

An individual computer system can be an NSW T8H and at the same time directly service non-NSW users. However, the resources supporting the T8H implementation are dedicated exclusively for this purpose, and cannot ordinarily be manipulated from outside the NSW environment. This helps maintain the integrity of the globally maintained system resource data base.

The administrative centralization of network resources in NSW means that procedures for establishing access rights to these distributed resources are quite simple. One merely negotiates with a single NSW administrative organization to arrange an account which potentially allows access to all NSW resources. The logical centralization of authentication, resource allocation and access control functions within a single component simplifies the Software implementation of these functions. It also means that there is a convenient framework for accounting for NSW resources utilized by an individual user.

The NSW attempts to shield the user from elmost ell details network operation, and emphatically tends toward the "invisible distribution" and of the design speatrum. This position is based on the feeling that dealing with the network or the constituent Host systems at any level would impair a user a ability to easily utilize combinations of moftware development tools. However, there are a number of factors which suggest that a more "visible" approach to network systems is sometimes more appropriate. For example, there is frequently an apprec(able decrease in performance when remote resources are accessed. Also, verious similar system resources may have slightly different characteristics on different Host types. still too early in the development of network=based distributed systems to essert with confidence that one can identify the best operating point in this (visible/(nvisible) design spectrum,

3.5.3 The TIP Service Fedility

As noted in the discussion of the TIP system (n Section 3,2, the users (and administrators) of a small computer such as the TIP will almost always desire more services than the small computer can provide. In particular, because of memory

limitations, the TIP is incapable of providing its users with a sophiaticated command language. The TIP has no space to hold tables of passwords or statistics on its usage; thus, the TIP has no dapability for access control or accounting. The TIP cannot distribute operational information to its users, such as announcements of system changes. What the TIP does provide is a relatively transparent, simple, flexible, and high performance interface between a terminal and the network. However, if access control, accounting, and operational dapabilities are to be provided, it is necessary to devise a mechanism to obtain these capabilities elsewhere.

During the development of the RSEXEC system two relevant observations were made. First, since many of the features planned for the RSEXEC were well matched to the desires of TIP users, with some additional effort the RSEXEC system could provide TIP users with a sophisticated sommand language and other features they desired (Cosell75). Second, because the RSEXEC was to be run on several systems, RSEXEC could potentially provide capabilities to the TIP very reliably. With a single Host providing a function, there would be times at which that Host would be down when some TIP user required the function. Thus, it would be possible through TIP use of the RSEXEC to obtain TIP

capabilities superior to any the TIP sould provide itself or that could be provided with the help of any single other Host.

Two mechanisms were developed to support the redundant implementation. The first is a "broadcast" Initial Connection Protocol. This enables a TIP to connect to an evaluable and responsive RSEXEC rather than to a particular version at a specific site. Using this medhanism, a TIP broadcasts requests for service to the known RSEXEC sites and then selects the site that responds first as the one to provide the service.

The second is a mechanism to meintain multiple copies of the various information files (e.g., network news and Host schedules) at the RSEXEC sites. This machanism allows additions to these distributed information files to be initiated from any RSEXEC site and guarantees that the additions are incorporated into each file image in a consistent manner.

Both of these mechanisms have been in operation for several years, supplying desired facilities to TIP users upon demand, Perhaps more interesting, however, was the extension of this system to several months of operational use in providing TIP access control and collecting TIP accounting information. These uses are currently dispontinged because of a number of broader

administrative issues, but while active the access control and accounting procedures function with a high degree of success.

During times when the mechanisms are active, the TIP code is changed so that whenever a user activates a TIP port, the TIP uses the broadcast ICP mechanism to connect to an R8EXEC which acts as a network login server. If the user successfully supplies a valid name and password, he is granted continued access to the TIP, the network, and to the standard RSEXEC functions. In addition, the R8EXEC transmits the user's network ID code (which serves to uniquely identify the user for authentication and subsequent accounting purposes) to the TIP and makes a "login" entry into a "TIP accounting" data file. If the user fails to supply a valid name and password then the R8EXEC does not forward an ID code to the TIP. If the TIP does not receive an ID code within the allowed time the TIP disconnects the terminal and the user is denied further access.

After the TIP redelves the user's network ID code it activates "connect time" and (outgoing) message counters to eccumulate usage data for the user's session. These counters remain active until the user terminates his TIP session. Periodically the TIP executes an "accounting checkpoint" procedure whereby it transmits usage data for its active users,

accumulated since the last checkpoint, to a data collection server process. The data collection server stores the checkpoint data in an incremental TIP accounting file for later processing.

Like the login servers, the data collection servers are redundantly implemented to insure high availability and to achieve load shering. The TIP uses a request mechanism similar the broadcast ICP to select one of the servers to accept its checkpoint data. The protocol is designed to allow considerable flexibility in the choice of a server. For example, a TIP can switch from one data collection server to enother after initially choosing one in the event that the chosen server can not complete the transaction (for example, because of network or Host failure).

The collection of ingremental accounting files created by the data collection servers is a large distributed and segmented data base. Some checkpoint data for a given TIP session may be collected by a data collection process at one site while other parts of the data for the same session are collected at other sites. The reduction of data in this distributed data base to produce periodic accounting summaries is accomplished by software which executes within the environment provided by the RSEXEC distributed file system.

The TIP service facility uses two fundamentally different types of distributed data bases. The first is one which is maintained "identically" at a number of sites. The second type consists of distributed, noneoverlapping segments; that is, the data base is a collection of segments, each of which is singly maintained at a (possibly) different location. These two types represent extremes of distributed data base organization; other applications may call for "intermediate" types such as a data base consisting of a collection of segments of which some, but not all, are redundantly maintained.

Use of the first type of distributed data base is exemplified by two aspects of the TIP service facility:

- The RSEXEC maintains a copy of the TIP news file at each of the RSEXEC sites. Updates to the news file are limited to addition of news items. The system allows additions to the data base to be initiated at any RSEXEC site and insures that all such updates are transmitted to and incorporated into all copies of the data base.
- 2. The TIP login system requires that the network user ID data base be maintained in a consistent menner at all RSEXEC sites. Each copy of this data base is a collection of

mutually independent user entries. Allowable updates to this data base include the addition, modification, and removal of individual user entries. The data base management techniques allow updates to be initiated at any site and guarantee that they are consistently incorporated into all copies of the data base (i.e., if all updating activity were to cease, all copies of the data base would aventually be identical).

The techniques used to maintain the data bases consist of two independent parts:

- A reliable, data-independent, update transmission and distribution mechanism which uses persistent processes at the update entry sites to guarantee that all updates are eventually delivered (once, and only once) to all data base sites.
- 2. A data-dependent update action procedure which is activated at data base sites whenever update commands arrive.

The operation of the TIP accounting system results in the creation and manipulation of the second type (i.e., segmented) of cats base. The primary sondern in the accounting application was

data base issues that required attention were:

- La Cataloging: It is obviously important to know where the various data segments (incremental accounting files) reside so that they can be accessed. This cataloging function is provided by the RSEXEC distributed file system.
- 2. Insuring that no duplicate entries occur in the data base, Because the entries contain accounting information, it is oritical that any redundancy does not cause duplicate charging. The data collection protocol was carefully designed to prevent the occurrence of duplicate data entries in apite of the fact that data is broadquat to all of the servers.
- 3. Insuring that each data base entry is processed exactly once when accounting summaries are produced. Time stamping plays a fundamental role in guaranteeing "once only" processing.

The implementation of the TIP service facility (a interesting primarily because it provides a working demonstration that it is possible to provide easy access to large and complex functions in a transparent way from much simpler mechines.

Further, the presence of multiple components in a distributed system, together with the potential for their redundancy, makes it possible to achieve reliability by constructing systems from modules most of which are kept relatively simple. By using simple modules, component failure due to melfunction of non-easential features can be reduced. Complex components are redundantly supported in an effort to enhance their reliability. Such an approach takes full advantage of both the heterogeneity and homogeneity of various network domponents. The important issues in designing a system of this type are the assignment of functions among the various mechines, the degree of redundancy required, and the protocols used to bind the system modules together.

Once such a virtual executive is conveniently available to TIP users, it becomes possible to think of additional features that can be added. For instance, the RSEXEC makes available to Host users a file system which spans machine boundaries. It is a simple technical step to provide the TIP users (who, unlike the Host users, have never had a file system) with a virtual file system. Another examples while the RSEXEC has the capability to permit users to leave messages for other users, it does not provide the capability for TIP users to receive such ressages.

Yet, through the concept of resource sharing, the potential capability to provide virtual mailboxes through which users can receive messages exists. Furthermore, through the redundancy inherent in the system, these mailboxes could be provided in a way which would insure that a user's mail was accessible no matter which individual computers were down. A final examples once the TIP user is connected to the RSEXEC and is ready to use the services of some Host, and once it is possible for the user to call for service independent of Host, there is no need to retain in the user's view the concept of the Host(s) from which service is obtained; rather, the virtual executive could be expanded to provide the virtual operating system from which all service is obtained. Such a concept combines the physical "front-end" properties of the TIP with the logical FE/MM partnership of the NSW.

3.6 Network-based Experiments and Research

The ARPANET subnet of IMPs and communications links was created primerily to provide a communication service to existing Host organizations and user populations. These groups have, over the period of network operation, been involved in many areas of research involving programs implemented on the Hosts; these programs may be categorized as "Hostebased" research. However, snother class of research projects and experiments have come about as a result of the existence of the ARPANET; this class may be characterized as "network=based" in that it seeks to evaluate or extend the utility of the packet switching concepts embodied in the ARPANET.

This section discusses several of the network-based experiments and research programs which were (netituted prior to the time the management of the network was assumed by DCA; many of these research programs are continuing to the present, It is also worth noting that an extensive program of research in extending packet techniques to radio networks is underway [Kahn75], but is not discussed here because its interaction with the ARPANET is minimal.

3.6.1 AFCS Experiment

Although ARPA itself provided funding of the Network Control Center, which measured network performance from an operational point of view, and the Network Measurement Center, which can enumber of performance experiments involving very large measurements, the first (and perhaps only) outside evaluative measurements of the ARPANET were carried out by the Air Force Communications Service in 1972 [AFCS72], The experiment was designed to evaluate the suitability of the ARPANET for transporting large volumes of traffic between pairs of sites; specifically, the performance and cost were to be compared to the Circuit Switching Unit (CSU) of AUTODIN.

The experiment was defined, in early 1971, to be based on the interconnection of existing Air Force Universal computers at Tinker Air Force Base (Okiahoma) and McClellan Air Force Base (California). ARPA agreed to install IMPs at these two locations, and assisted the AFCS in the definition of an ad hoc Host-Host protocol for use in the experiment. The protocol eventually defined used multiple link numbers for a single communication in order to achieve high bandwidth, and used a minimal Host-Host header in order to achieve high efficiency.

The actual measurements took place from mid-April to the and of June in 1972. The shortest path between the Tinker and McClellan IMPs during this period traversed five hops (i.e., five IMP=to=IMP eircuits.

The test findings and conclusions which resulted from this experiment were:

- 1. Throughput between the two Hosts for all traffic averaged 25 Kbs.
- 2. For test treffic of known contents, 31,761 18mbit words out of 247,365,000 words transmitted were observed to be in error; this amounts to about 01%. Analysis led to the conclusion that most (or all) of these errors were introduced in the IMP/Host interface. Therefore, AFCS recommended that a hardware error detection scheme be implemented for this interface: In addition, it was recommended that a software error detection scheme be incorporated in the Host-Host protocol:
- 3. The coat per megabit transmitted was found to be approximately equal to that of the AUTODIN CSU at daily traffic levels below 168 megabits. At higher daily

traffic levels the ARPANET was less coatly than AUTODIN,

4. It was recommended that further studies be conducted "to determine the feesibility of implementing a secure network, patterned after the ARPANET, for DoD users with large digital transmission requirements" [AFC572].

3.6.2 Speech Transmission

There are several reasons for interest in transmitting speech signals in digital form through a pecket-switched network. Perhaps the most important stems from an intrinsic difference between traditional circuit switched transmission and packet switched transmission as exemplified by the ARPANET; during periods of excessive demand circuit switched systems ber access to new calls, whereas packet switched systems may merely increase per message delay. The advantage of a guaranteed connection is so important in some government and commercial situations that work dedicated lines are lessed in order to provide it (Forgie?5). The advantage of a guarantee at lower cost.

In addition, digitized speech offers other advantages.
Digital transmission is less sensitive to noise, encoding a

digital signal for privacy or security is masier, and buffering, processing, and archiving spoken messages on large Host computers is possible.

A major obstacle to transmission of speech through packet networks is the number of bits per second required to encode the speech signal, High fidelity analogato-digital conversion would result in about 250 kbs; provision of telephone-quality speech signal would require about 50 Kbs. For this reason there is currently a great deal of research in the area of speech bandwidth compression. There are two general approaches under investigation; these can be denoted the waveform coding eppreach and the vocoder approach. The waveform coding approach is based on the fact that successive samples of the speech waveform are not independent; allowing a restriction of the range of possible successors to a given state which the coding must represent, and thus a restriction in the code size. The uncoder approach is based on methods of modeling the speech process. If the modeling is suggessful, a speech signal can be represented by a set of parameters which require lower bandwidth for transmission than the signal (tself, Demonstration vocoders have required as few on 1,2 Kbs for the signal characteristics of selected speakers, and adequate goding of most voices can be done In the 2.4=9.6 Kbs range.

ARPANET packet speech experiments have been based on the use of relatively low bandwidth, speech compression devices, mostly in the vocoder class. However, the majority of available devices have been constructed assuming connection to switched (or dedicated) communication channels of fixed bandwidth and delay; these characteristics, of course, do not pertain to the ARPANET. Thus the packet speech research and development has been designed to try to find solutions to two problems:

- 1. Finding ways to reduce the average transmission rate by taking advantage of the variability of the actual apeach data rate, in order to avoid unnecessary network load.
- 2. Finding weys to meak the veriability of ARPANET transmission delay when regenerating the speech signal; in order to evoid listener annoyance.

By increasing the complexity of the vocader's interface to the communication channel, in perticular by interfacing it to the ARPANET through a Host of modest size, it has been relatively easy to reduce the average data rate. For example, no bits at all need to be transmitted when the talker is silent, either pausing to think or to wait for the other party in a conversation to finish talking. Further reductions in average data rate are

possible by taking advantage of the fact that during certain speech sound the character of the sound thanges much more slowly than its maximum rate.

It has been less sesy to mask the variability of ARPANET transmission delay. Obviously, some buffering is required at each vocaders at the source buffering in required to construct packers and hold them until the local IMP accepts them, and at the destination buffering is required to read the packets and feed the bits serially to the speech synthesizer, appeared that by expanding the deskination buffering to handle seconds of signs) (a few thousand bytes of storage), and starting the synthesizer only when the buffer was filled (e.g., half) with input from the network, the veriability could mostly be masked from the listener. However, the IMP to destination IMP algorithms were designed for very reliable transmission; to aphieve the desired level of reliability packets may be retrensmitted many times in order to be sure they get through. In addition, the ARPANET promises to deliver data in the same order in which it was transmitted. Thus, if a packet must be retransmitted, the delivery of all subsequent packets will be delayed. Since source IMP to destination IMP (noutry and retransmission timeouts may be measured in tens of seconds, a

this is still a good technique with 8 is the proper Forgie 50/20 Hunks

mingle last packet may have a major disruptive effect on the appech receiver.

A human listener, of course, is well adapted to interpolating speech signal across bursts of noise. Thus, from the point of view of packet speech transmission, it would be preferable if the ARPANET did not try so hard to deliver the date reliably. An algorithm which gave up after a second or two and delivered subsequent packets without worrying about the gap would be much better. The protocol used for the transmission of speech packets could include timing information, and the synthesis device could fill signal gaps with a continuation of the immediately-preceding signal, white noise, or complete silence.

To summarize, the packet speech transmission requirements are neither for high bandwidth, low delay, nor almost perfect reliability; rather the requirements are for medium bandwidth, modest reliability, but very low veriability of delay. Special HosteHost protocols can be used to reorder packets; detect missing packets, smooth the delay, and adjust the output signal appropriately. In fact, the ARPANET was modified to allow a restricted set of Hosts to send packet traffic which bypasses the normal source IMP to destination IMP protocols, precisely for

peaket speech experiments. Unfortunately, however, the IMP congestion control algorithms are intimately entwined with the missing packet detection and retransmission algorithms; thus unrestricted growth in the speech traffic might lead to congestion problems or other interference with the rest of the Hasts. Nevertheless, several demonstration conversations have been carried out over the ARPANET using these techniques.

3.6.3 Packet Broadcast by Satelifte

The ALOHA System (Abramson70], developed at the University of Haweii, introduced the use of a "multimecoess" channel for packet communication. In this system, a number of independent transmitters accumulate packets of information; all packets are to be sent to a single central site, Whenever a transmitter has accumulated a full packet it appends its own address and a checksum, and transmits this packet on a fixed radio frequency. The same frequency is used by all transmitters, so two transmissions may overlap in time, If this happens, the checksum received at the central site will indicate an error and the (composite) transmission will be (gnored.

A second frequency is used for all packet transmissions from the central site to the remote stations. Since this frequency is

used by only one transmitter, conflicts do not occur. Address information in these packets allows the remote receivers to identify the packets intended for them. Packets correctly received at the central site are answered with positive acknowledgments sent along with central site data. Thus, when a remote station sends its packet it starts a timer; if the timer expires before the acknowledgment arrives the packet must have been destroyed by channel noise or conflict and is retransmitted.

The capacity of the multi-access channel is obviously less than the rated channel bandwidth, due to conflicts and retransmissions. In fact, under certain reasonable assumptions about frequency distribution of pecket errivels, it has been shown that a channel operated in this way has a capacity of about 18% of rated bandwidth. However, it has been noted that if such an ALOHA channel were divided into slots (each able to hold a packet) and if the nodes modified their behavior to transmit packets so that the leading edge of the packet always coincides with the leading edge of a slot, even though the nodes remain free to transmit into a slot without regard for the transmission of other nodes, the effective channel capacity is doubled to 36%; A channel operated in this manner is called a "aletted ALOHA" channel.

Packet broadcast by satellite extends these ALOHA concepts to a geosynchronous satellite transponder and an ARPANET-like environment. The principal differences are:

- 1. There is no central site. All stations are equal in their access to the channel; each accepts only the packets addressed to it.
- 2. The long roundatrip time makes it worthwhile to ettempt to schedule portions of the channel (in a highly dynamic way) for any port(one of the traff(c with known characteristies.
- 3. Since there is no conflict-free frequency portion of the channel on which to send positive asknowledgments, it may be worthwhile to achedule a gonflict=free time portion of the ghannel for each transmitter to gend acknowledgments. This may conserve bandwidth by helping prevent acknowledgments from being lost, with needless retransmissions resulting from such loss,

In line with theme ideas, ARPA has commissioned experimental packet matellite network to darry out tests of a variety of channel protocols and operational nodes. This network required the development of a new Satellite IMP which controls the use of the channel. The Satellite IMP is a conventional IMP with several additions and modifications both in hardware and software (Butterfield74).

The first hardware addition is to increase the memory to a size sufficient to buffer all the transmitted packets which can be awaiting adknowledgment simultaneously. Suffer space is necessary for some 32 packets assuming a 50 Kbs channel and a one quarter-second propagation up to the synchronous satellite and back down. That is, 32 packets can be sent out before the acknowledgment returns for the first.

The next herdware addition is a mechanism for signaling the satellite radio transmitter when to turn the radio carrier on and off as packets are transmitted; this is necessary because if two satellite ground stations have their radio transmitter carriers on simultaneously, they law each other.

The third addition is time-keeping hardware necessary for the Satellite IMPs to accomplish accurate slatting. This hardware notes the arrival time of the leading edge of a packet (slot) and makes this time available to the program when the program fields the received packet interrupt and updates the

program's estimate of the slot positions. This hardware also allows the program to accurately specify transmission of a packet at a specified time in the future.

One hardware modification was nacessary for construction of the Satellite IMP. The IMP modem interface normally uses a unique character sequence to denote the end of a packet, thus requiring an escape character with escape character doubling for data transparency. Escape character doubling, however, can result in the length of a packet being temporarily increased while traversing the satellite, thus overflowing a slot. The IMP modem interface, therefore, had to be modified to use a word count to specify packet length.

The Satellite IMP software effectively makes the Satellite IMP an extension of an ARPANET IMP, using the satellite channel to provide an additional ARPANET link between the Satellite IMPs, The satellite channel protocol, however, is designed to allow many Satellite IMPs to share the channel. Because of its special nature, the satellite channel and its associated Satellite IMPs are most conveniently treated as an independent network temporarily interconnected with the ARPANET.

For experimental purposes, one of four different multiple access broadcast channel protocols can be selected by a software switch in the Satellite IMPs. These protocols consist of fixed TDMA, slotted Aloha, a rudimentary version of Reservation-Aloha, and Reservation-TDMA. Each of these is a demand assignment (DA) system, although the fixed TDMA protocol allows the DA only through its varieble destination (packet address) capability. The other three also have this capability, but in addition allow channel transmission time to be shared among the Satellite IMPs in a completely dynamic manner,

Sased on some experience with these protocols, as well as continuing analytic work, a contention-based priority-oriented demand assignment, or CPODA, channel protocol has been under development. The goal of the CPODA protocol is to provide a very flexible multiple-access demand-assignment system for use in a broadcast satellite channel, while also allowing users to specify a number of different priority levels and delay constraints, further, the protocol is designed to allow the integration of both block and stream traffic, where the latter refers to packet voice and similar types of traffic having a relatively short delay constraint on each packet along with a relatively constant bandwidth requirement over the time span of a conversation.

A major feature of the CPODA protocol, which distinguishes it from the Reservation=TDMA protocol currently implemented in the Satellite IMP, is its use of a rendom access technique for making reservations. Channel time is divided into subframes, with one subframe used in a rendom access manner by all stations to make reservations, while the other subframe is used for the resulting acheduled transmissions of both block and stream data packets.

To allow experiments to be performed with these protocols unencumbered by ARPANET source=destination protocols and uncontrolled traffic sources, the Satellite IMP software also includes the following additions:

- Satellite IMP Message Generators This code provides a controlled source of experiment traffic which is not subject to IMP source-destination protocol, and provides a wide range of message lengths and rates.
- One-way Dalay Measurement Code: This code makes use of the synchronized clocks provided by the Satellite IMP's slotting algorithm to measure the delay encountered by each measage between its time of origin in one Satellite IMP and its correct receipt in the destination Satellite IMP.

Access Control Code: This code overrides normal IMP routing decisions, allowing only experimental traffic to use the satellite channel while forcing all other traffic over normal ARPANET circuits. To allow use of the satellite channel as a backup link for ARPANET when the normal transmatlantic path is disrupted, the code automatically switches, when necessary, into a "normal" routing mode.

The Satellite network consists of two Satellite IMPs located at large-entenna (30 meter) earth stations in Etam, West Virginia and Goonhilly, England which have been in place for some time, The system uses a 50 Kbs satellite channel supplied through the SPADE subsystem. Two additional nodes were installed in the Satellite Network later in the project; these are located at the Scandinavian large-antenna site in Tenum, Sweden and at a COMSAT laboratory small-entenna site in Clarksburg, Maryland. The Etam, Goonhilly, and Jenum sites are each directly connected to the ARPANET via IMP+to=IMP circuits, so for a considerable period of time these Satellite IMps have been operated as an extension (or subnet) of the ARPANET proper. However, the long-term intent is to operate the Satellite Network as an independent, although interconnected, system.

3.6.4 Internetting and Gateways

The work done to date on the interconnection of computer natworks has usually assumed a configuration such as that shown in Figure 3-3. On each natwork there are Hosta (denoted by H in the figure) which desire to communicate with Hosta on other natworks. The natworks are connected together by units (denoted by G) called "gateways". The gateways must in some way convert traffic in the format of one natwork into traffic in the format of another network.

Because Most-Most protogols differ from one network to the next, and because these protogols are generally complicated and incompatible, it seems clear that Hosts on different networks wishing to communicate must do so in a common protogols. Therefore, much of the work to date in network interconnection has been the specifications within such bodies as IFIP Working Group 6:1 [Cerf76] and more recently ISO of dommon and-to-and protogols, However, since standardization efforts tend to take a long time, ARPANET experimenters have "standardized" on an end-to-and protogol described in 1974 [Cerf74]. In this protogol the logical antity in the Most which performs the protogol functions is called the Transmission Control Program or TCP.

Figure 3-31 Networks Connected by Gataways

Given a Host protocol (TCP) and a pair of networks of somewhat different characteristics (ARPANET and the Satellite

Network described in the preceding section), ARPA has been involved in the definition of experiments to investigate the functions that gateways should perform in a network interconnection environment. In fact, by the time that the ARPANET to Satellite Network gateways (three of them) achieved initial operation in mid=1977, plans had been formalized for pairwise gateway betweens

- ARPANET and ARPA's experimental packet radio network
- ARPANET and an internal network at 88Nfs Research
 Computer Center
- ARPA's experimental packet radio network and ETHERNET at the Palo Alto Research Center of Xerox
- ARPANET and an internal network at MIT's Laboratory for Computer Science

as well as probably expansion of the ARPANET to Satellite Natwork gateway in England by addition of a connection to the Experimental Packet Switched Service of the British Post Office.

One of the outstanding questions of network interconnection is whether the gateways should connect networks at the packet or Host level. At packet level, a portion of the gateway would actually become a node on each of the networks being connected,

while at Host level, a portion of the gateway would actually be a Host on each of the networks being connected. In the current experiments, the gateways connect at the Host level primarily to maintain a "sovereignty" of the networks involved. This choice allows connection to the networks at a point where the interface is both well defined and well controlled by each constituent network. Each network can protect itself against activities of the gateway to the same extent as it may protect itself against the activities of the activities of any other Host.

Since the gateways connect to the networks as Hosts, then the format of the messages passed to a network are specified by the Host/network protocol. This protocol is then used to permit transparent transmission of segments of an internetwork message by embedding the internetwork segment in the text of a local network message. Such a composite unit has two leaders and potentially two trailers. The outermost leader and trailer provide information for the local network and will be different in each network through which the message passes. The leader will specify the address of the gateway Host to which the message should be delivered, any allocation or sequencing information which is used by the Host/network protocol, and any further information demanded by that protocol. An example of a trailer

Chapter III

that might be required by the Hast/network protocol would be padding and a checksum. Within this outermost leader and trailer is the internetwork data segment with its leader and trailer, The internetwork leader specifies such information as the ultimate destination, sequencing, and reassembly information, The actual data which is being transferred is the text of this segment.

In addition to the ability to accept messages from one network and resend them into another, it is likely to be desirable for the gateways to have the following additional characteriatical

Routing Capability: Inter-gateway routing is desirable for all of the standard reasons routing is desirable in a network, For example, for reliability there abould be alternate paths over which traffic may be routed; for achieving higher bandwidth than is available over any single path, gateways should be able to route traffic over parallel different classes of traffic need to follow different routes (e.g., traffic requiring low delay should be routed around networks which insert large delays).

- Access Control and Accounting Mechanisms: A given constituent network may wish to limit some classes of traffic or all traffic at some times (e.g., because of regulatory considerations, dountry A might not want traffic from country B to dountry C to pass through dountry Afs network), Also, a given user might not wish his traffic to pass through some particular constituent networks, e.g., for political reasons. Alternatively, the constituent networks are very likely to want to monitor the use of their network by internations traffic.
- Ability to Perform Message Fragmentation: Because of the differences in message size of the constituent networks connected by a gateway, the gateway must have the ability to fragment a larger message arriving from one network into smaller messages which are acceptable by the next network. When such fragmentation occurs, the message stream must eventually be reassembled into its original structure. The TCP protocol provides this function at the destination Host.
- Congestion Control: Congestion will inevitably occur at the gateway unless specific measures are taken to prevent it;

 This congestion can occur as a result of speed mismatches

between the networks connected by a gateway, because several gateways on a network may simultaneously transmit traffic to a single destination gateway, because traffic may have to be held during a period of recovery from a failure, and so on. One specific kind of congestion results from deadlooks, such as when gateway A is full of traffic for gateway B which is full of traffic for gateway A.

Retransmission Capability: When a message is lost in the network between two gateways, either the last gateway can retransmit the message or the message can be retransmitted from the source Host to the destination Host. It has been shown that hopetoehop retransmission is more efficient than source=to=destination retrensmission if the possibility of message loss is appreciable; and even when there is little possibility of message loss, the variance of retransmission delays is less with hopetowhop retransmission than with sourcesto-destination retransmission. While some networks deliver messages very reliably, other networks rely on source to destination retransmission and in some cases are quite cavalier about throwing away messages. Thus, it seems that at least when the Hosts on a network are normally responsible for retransmission across that network, the gateways ought to provide retransmission.

A gateway with the characteristics described above is very similar to a node on a packet switching network. This leads to the notion of a gateway virtual network wherein the gateways act as nodes and the network spanned by the gateways acts as virtual lines fully connecting all the gateways on that network. Further, there is a gateway logically associated with each Host attempting internetwork communication, providing at least some of the functions (e.g., routing, but perhaps not retransmission) which must be provided between networks. However, there need not necessarily be a one-to-one correspondence between Hosts gateway machines. For instance, the logical entity that is the gateway may have the form of a program running in a Host computer. Alternately, the gateway could be in a standwalone machine serving one or more Hosts. On the other hand the gateway connecting two networks could take the form of a program running in a Host which is connected to both networks. In General, a dateway should be able to connect any combination of Hosts and networks,

It is entirely possible for a Host not to have its own gateway, preferring to use a gateway elsewhere in its network to perform the gateway functions it needs. In this case, the Host would simply know the addresses of a few gateways in its network,

and would send its internetwork traffic arbitrarily to one of these gateways for routing and forwarding. Of course, the gateway arbitrarily chosen might not be on the best path to the destination, and this gateway might even forward the traffic to another (better) gateway in the same network for forwarding outside the network. Thus, this approach may have some inefficiencies but it reduces the number of gateways that have to be constructed.

The gateway network has many of the maintenance characteristics of a standwalone packet switching network. It requires centralized development and maintenance responsibility, including a gateway network monitoring and control center.

The Gateway Virtual Network solution is a general solution to the problem of interconnecting networks which is not highly dependent on the nature of the networks being connected. Because of this it can be expected that additional networks which are connected in the future will not require modification. A principal goal of the Gateway and Internetting experiments currently underway is to discover whether this expectation can be met.

3,6,5 End-to-End Security

As described in Section 3.3, it has been possible to build equipment (the PLIs) which allows use of the ARPANET for the transmission of Classified data in a secure way between pairs of Hosts, and it appears likely that use of a multi-address capability for this equipment will be obtained.

Nevertheless, the PLI (even the multi-address version) is far from ideally matched to the possibilities offered by the network. For example, the PLI is too large and expensive to be useful for terminal connections. The KG units at all PLIs must use a single key to allow general communication, however, it might be desirable for A and B to engage in communication secure from C and for A and C to engage in communication secure from B.

for reasons such as these ARPA has engaged in a program of research on end-toward security, in an unsecured packet network any (ronment, with the following goals:

- Provide secure communications over a packetmoriented network
- Provide end-to-end encryption for both Host computers and terminal users

- Avoid the use of classified equipment in the experimental ayatem
- Provide a system as close to militarily secure as possible, with an eye to certification of a similar system
- Achieve efficient use of the packstmoriented network
- Design a system which could be deployed cost effectively
- Operate in an internationking anvironment, lliustrated with at least the ARPANET, the Packet Radio Network, and the ARPA Satellite Network
- Make use of proven techniques where possible
- Minimize impact on Most operating systems; and, in fact, if possible make the security system transparent from the Hosts? point of view.

The encryption unit developed in this program is referred to as a 8CR. The letter 8 indicates the black side of the unit, the letter C an interior crypto or KG portion of the unit, and the letter R the red side of the unit. This three component unit taken as a whole provides a secure interface between the Hoet and the network. The B portion of the unit is involved both in interfacing between the BCR and the IMP and in an endeto-end (8 to 8) protocol. The R portion of the unit is involved both in

interfacing between the BCR and the Host (or term(na)) and in an end-to-end (R to R) protocol.

Each BCR contains two LSI=11s, one as the Black side processor and one as the Red side Processor. The Crypto unit is a fairly complicated unit with storage for several different Keys, control paths which allow Keys to be loaded and selected, control paths which couple the Red and Black processors with low-bandwidth unencrypted paths, and a KG unit. The actual encryption algorithm chosen for the KG unit is the algorithm now adopted as a federal standard. The BCR data and control flow paths are shown in Figure 3-4.

Figure 3-5 is useful in describing the general approach of the experiment. There are four different site configurations illustrated from top to bottom. At the top is shown a typical secure Host. This Host supports general processes. The processes interface to the network through a user level protocol such as TELNET, and the TCP (see the previous section). Host-Host protocol. The Host itself is interfaced to the packet network via the BCR. The second configuration from the top of the figure is very similar to the general Host configuration except that the Host, instead of supporting general processes, supports a number

DRAFT Expertence

Figure 3-4: BCR Data and Control Patha

DRAFT Experience

Figure 3+5: Possible Site Configurations

of local terminals. Two Hosts configured as in the top of the figure, or two terminals concentrators configured as in the second line of the figure, or a terminal concentrator and a queeral Host, can communicate in a secure manner if their respective KG units are using matching keys. Furthermore, a given general Host or terminal concentrator can simultaneously communicate securely with more than one other such site, since the BCR units are able to select a different key for each message by a method described below.

The third configuration illustrates a situation where a terminal is not so near other terminals that it is don't effective to install a terminal concentrator Host; this situation is indicated by the terminal connection to the packet network being via modems. In this case, the BCR unit and the TCP and TELNET functions must all be effectively at the terminal. This currently requires three LSI=11s; the TCP and TELNET are incorporated into one and the others serve as 8 and 8 processors.

At the bottom of the figure (a illustrated a special Host called a Key Distribution Center (KDC). This Host takes part in setting up the correct keys in the verious BCR units. To be able to properly set up the keys, the KDC must be able to

simultaneously communicate with the red mide of the BCR units and the black side of the BCR units. Therefore, the KDC must simultaneously be connected to the network both with a BCR (for communication with Rfs) and without a BCR (for communication with Bfs). This dual connection has the further implication that the KDC itself must be certified to not look secure data. If a communicating Host is sufficiently cooperative, the KDC can provide very specific security control, down to the level of security compartment, need to know, etc. The KDC for the current experiments is implemented on a TENEX system.

A scenario of the operation of the system is as follows, When a data packet is presented to the red minicomputer for transmission over the metwork, it must include (or imply) its source, destination, and length. Furthermore, if the packet is the first such packet, information for determining the classification of the data must be included. If the red minicomputer discovers that there is no entry in its tables relating the specified virtual source and destination to a link (described below) it either discards the packet (if no classification information is supplied) or sends a message to a KOC, using a preset link through the KG unit.

In response, the KOC sets up a duplex communications channel from end to and through the network. This is done by instructing the black minicomputer to set up the key storage for an unused link in the KG units at each and of the secure channel and specifying the actual source and destination addresses to be associated with that link in the black minicomputers. It also informs the red minicomputers of the link to use for that virtual connection. Now packets can be sent over the link until either the red minicomputer (from information supplied in the packets) or the KDC closes the link.

Figure 3-6 shows the relationship of the addressing tables loaded into the minicomputers and the keys loaded into the KG units to the link number. The common piece of information is the link. The link can be used to select the encryption key, the decryption key (they need not be the same), the virtual address information, and the actual address information. The address tables in the minicomputers permit the inverse translation to be done as well. One or more links are predefined to permit messages to be sent to the KDC.

The operation of the red and black minigomputers is nearly symmetric for outgoing and incoming messages. When the red

Figure 3-6: Relationship Between Keya, Addresses, and Links

(black) minicomputer receives a packet from the Hoat (network) (tatrips off the virtual (actual) address information and looks up the link to use for that connection. It supplies that link number to the KG unit and then passes the remainder of the packet into the KG for encryption (decryption). The KG passes the link number on to the black (red) minicomputer and then the encrypted (decrypted) data. The black (red) minicomputer uses the link supplied to insert the actual (virtual) address information into the packet and sends the packet to the network (Hest).

The differences in operation of the red and black minicomputers are in the handling of packets for which no corresponding link can be found. The red minicomputer attempts to establish the link as described serlies but the black minicomputer discards the packet.

The modification to the Host-Host protocols necessary to use this secured network are minimal. Principally, the source TCP must supply classification information to the red minicomputer to permit it to establish the necessary link. This information is encrypted and passed along to the destination where it may be used as an additional check on the validity of the connection. This information may be supplied at any time but must be supplied

Section 3

DRAFT Experience

in the initial packet. In addition, the red minicomputer at one end of the connection (at least) must be informed by its Host computer when the link is no longer needed. This information is not the same as the information which the two Host computers exchange to close the connection, since packets containing that information may need to be retransmitted.

3.7 ARPANET Operation and Maintenance

The operation and maintenance of the ARPANET (tself makes extensive use of the communication subnet and a number of network Hosts. The highly centralized approaches taken to these functions are quite novel, and a number of sophisticated, even if ad hoc, techniques have been developed during the course of operational experience. However, unlike most of the other significant aspects of network construction, experimentation, and utilization, the network operating mechanisms and techniques were not designed from scretch, but rather grew in response to various internal and external pressures. Thus a discussion of network operation and maintenance would be incomplete without mention of some of these historical factors.

3.7.1 Control Functions in the IMP

The recognition of the necessity for network control functions began with the design of the IMP hardware and software in 1969 [McKenzie72] [Mckenzie75]. In the hardware design (t was recognized that with several parties responsible for vertous places of equipment in the network, the problem of fault isolation was critical. Accordingly, both the Host interfaces and the modem interfaces included the capability for "loopback"

test mode under progrem control. The loopback mode of operation basically connects the IMP's outputs to the IMP's (nputs so that the IMP may generate test traffic, send it through the interface, and shack the returning traffic against the traffic which was generated. Thus, for axample, if a carrier=supplied communications circuit connecting two IMPs appears to be giving trouble, the IMPs can be directed to loop both of the modem interfaces connected to this circuit. If test traffic passes through these looped interfaces successfully, then it reasonably certain that the carrier's equipment is at fault. 0 n the other hand, if one of the two looped interfaces continues show a pattern of trouble, it is reasonably certain that the IMP hardware or software is at fault. In either case, the appropriate repair and maintenence organization can be contacted, with little risk that the trouble will eventually be traced to equipment which is the responsibility of another party.

The IMP also includes hardware for the automatic detention of power failures and for automatically restarting the program when power is eventually returned. The IMP hardware (in the initial Si6 systems) also contained a "watchdog timer" which counted down a clock at a rather slow rate and, if the clock wesever counted to zero, generated a very high priority interrupt.

The IMP software periodically reset this clock to (ta maximum value when the program was operating correctly. The interrupt routine triggered by the watchdog timer automatically reloaded (from a neighbor) and restarted the IMP program. In this way it was possible to recover from many erroneous states, whether caused by hardware failure or software bugs, without manual intervention.

The original IMP software, in addition to the watchdog timer routine and the softwere necessary to loop and unloop interfaces and to send test traffig and check the results, agentained several other diagnostic features. These fall under two main headings, IMP statistics and ODT. The IMP statistics package contains a variety of routines, most of which are usable by any Hosts these Include the ability to trace the progress of a pagket through the network, the accumulation of message and packet length histograms, the ability to take a "snapshot" of the internal state of an IMP at a given instant, and the ability to generate artificial traffic. One statistics routine, however, called the "trouble reports" routine, has its use restricted to network control and manitoring functions and unlike other statistics routines, cannot be activated and deactivated by any Rost, but permenently enabled and reports only to a single лема (ра

location. DDT is a small debugging package which contains a simple command interpreter, capable of such functions as examining and modifying a memory word, clearing a block of memory, searching memory for a particular stored value, etc. DDT is structured so that it can be driven remotely through the network, returning any responses back through the network. Controls (which were not present in the initial implementation) insure that DDT can be operated only from appropriate locations.

3.7.2 Early Operation

During the first several months of the ARPANET, when all the installed IMPs were located in California and Utah, operation of the network relied heavily on the essistance of personnel employed at these sites, as well as on the more or less continuel on-site presence of the development staff members from BBN, Recovery from an IMP failure generally involved local use of DDT through the IMP console Teletype by site personnel, working in conjunction with a BBN programmer (either on-site or by telephone). Each IMP was equipped with a high-speed paper tape reader and each site kept a paper tape of the most recent IMP software version. In case recovery was not possible using DDT, site personnel reloaded the IMP through the tape reader. If

after Several attempted loads the IMP continued to fall, the problem was presumed to be a hardware problem and Honeywell field Service, which had maintenance contracts on all the machines, was called in to perform appropriate diagnosis and repair.

Each IMP monitored the status of the common carrier circuits connected to it and displayed their status in a set of lights on the IMP front panel. Site personnel would periodically examine these lights and, if the circuit was seen to be in a down condition, would use DOT via the IMP console Teletype to loop the interfaces, generate test traffic, and monitor the results, which again were displayed as an up or down condition in the front panel lights. Site personnel would then contact BBN and, depending on the results of the testing, either Honeywell or the common carrier would be notified of an equipment feiture.

However, with the installation of the network's fifth IMP at Bolt Beranek and Newman in early 1970, testing and disgnosis became considerably more centralized. The "trouble reports" statistics program in each IMP was set to send periodic reports to the canadle Teletype connected to the 88N IMP. Data on IMP status and circuit status were sent as ASCII text to be typed out on the Teletype. Since there is limited buffering within the

functioning completely.

IMP, and the output was produced on a low-bandwidth device, space within the message was at a premium. However, a person was required to read and interpret the data; the massage format thus had to be at least somewhat intelligible. Balancing these two factors resulted in concentrating on only that data which was vital to the continued functioning of the networks circuits usable or unusable, the state of the donsole switches at each IMP, or (by inference from the absence of any report at all, for a long period of time) indication that an IMP had stopped

Members of the IMP herdward and software development team were responsible for scanning this typescript occasionally and initiating appropriate actions when necessary. Since no individual was exclusively assigned to monitor the Taletype, and since personnel were normally available only during regular working hours, it was possible to have a long period of outage for a circuit or an IMP without anyone noticing it or beginning corrective action. However, since the network was very small, and since its operation was not critical to almost any Host activity, this did not constitute a severe problem.

As the network became larger and more reliable, the proportion of status messages typed on the Teletype which said anything other than "everything still DK here" decreased, thus making the location of messages which required action more difficult. In order to make location of these critical messages easier, the format of the trouble report routine was changed so that it used the "synchronized clock" of the statistics package, Each IMP was programmed to send a status message every 15 minutes; each IMP also examined its own status every minute and sent an additional message if it detected a status change. Thus every 15 minutes a block of "checkein" reports was typed, one from each IMP, Interspersed between these blocks on the typescript there would be an occasional status change report. In general, only the change reports required action.

3,7,3 The Network Control Center

As the network continued to expend in size, scenning the typesoript on the IMP Teletype for network events which required human intervention became increasingly difficult. In addition, there was considerable interest in periodic reports on IMP and line performance, statistics on circuit usage in order to obtain advance warning of approaching saturation, and Host traffic in

order to investigate accounting algorithms for network usage, These factors led to the attachment of a mmall Host to the BBN IMP dedicated to monitoring network performance and doing much of the bookkeeping required for the reports, as well as watching for events which might require human intervention. The computer chosen for this Host, known as the Network Control Center (NCC) Host, were (dentical to a network IMP, namely a Honeywell 316 computer with Host and modem interfaces, 12K of 16-bit memory, and a real-time clock. Choosing this particular set of hardware meant that, in the event of a failure in the NCC Host, the prototype IMP normally used for software development could be temporarily substituted. A set of peripheral interfaces, which were packaged separately from the NCC Host (tself, were provided to drive output Teletypes, a set of display lights, and operator switches; this equipment could easily be connected to the I/O bus of an alternate CPU if the need arose,

By mid=1971 the NCC Host had become operational and it become possible to make several changes to the format of the "trouble reports" generated by each IMP. First, the ASCII text format was abandoned in favor of binary encoding. Next, the reports were expanded to include more internal status information, as well as statistics on Host traffic and line

usage, It was also possible to increase the frequency of reporting to once a minute for the "check=in" and to send change notices as soon as changes were detected. The NCC Host computer was given the Job of examining each of these reports, noticing changes in network status, and alerting the operations staff which by this time had grown to dedicated operator coverage for about 12 hours a day, five days a week. Operators were informed by the NCC Host that network status had changed by printeuts on one of the Teletypes connected to the Host and, in the case of IMP or circuit failure, by flashing a light (a separate light for each IMP or circuit) and sounding an audible alarm. The NCC Host also accumulated the line usage and Host traffic statistics and summerized these on a second Teletype.

It is worth noting that the term "Network Control Center" is a misnomer; "Network Coordination Center" or "Network Monitoring Center" would be a more appropriate name. In fact, a very important result of the ARPANET project is the demonstration that it is possible for a network to dontrol and operate (taelf without explicit control from a control center. Many network designs include at their heart the assumption that there must be a manned control center with the capability of adjusting the notwork's routing, instructing nodes which lines to use,

adjusting nodel buffer storage, adjusting nodel priorities, and so forth. The ARPANET depends on no such dentral gontrol. All such adaptations are made by the nodes themselves, alone or in concert; further, these dedisions are generally made more rapidly and accurately than could be done from a manned control center. The ARPANET NCG is responsible for coordinating maintenance and growth, but is not involved in real-time control.

functions of the Network Control Center Host have expanded rather dramatically, although in small steps, from the time of its initial implementation. First, the hardware was expanded from that of the standard IMP to that of the standard Terminal IMP, providing additional core memory as well as the possibility of operating several terminal output devices. late 1974 the Teletype which was used to print the log of network events was replaced by a 1200-baud line printer. This was necessitated not only by the greatly increased number of network nodes, but also by the large amount of additional information which the IMPs were reporting to the NCC Host. Reported events now include the status and error frequency of communications circuits in the network, the up/down transitions of each Host in the network, the lengths of various buffer queues (n each IMP, the settings of the various console switches on the IMPs, and a variety of "trap" conditions which may occur,

The trap mechanism allows the programmers or maintelners to have the IMP notify the Network Control Center whenever a particular block of code is executed. A single trap subroutine reports the program counter, the contents of the accumulator, and the contents of the IMP's index register for each trap condition. This mechanism is useful for determining the frequency of execution of various blocks of code within the IMP, and the set of traps changes as the attention of the programmers moves from one area of the IMP algorithm to another. The trap mechanism is also used to report failures of the software checksums and other difficulties related to hardware failures.

As the amount of information presented to the NGC operators on the log printer increased, and because an increasing amount of that information did not call for direct operator action at all, but was destined instead for system programmers and the maintenance staff, it became desirable to find a way to separate the log into several pieces. One possible solution would have been to connect a multiplicity of output devices to the NGC Host, each output device maintaining a separate log. However, the separation of the log output into several pieces need not be done in real time. Thus a single log is maintained, with the lights display and audible sparm as the primary method of

notifying the operators that immediate action may be required, Programs have been written for network TENEX systems which periodically probe the NCC Host and request copies of all available log information, which is then placed on TENEX bulk storage. Further, the data from the TENEX is sent to the Datacomputer roughly once per day for long-term storage. Programs have been written to examine this data and separate it into as many categories as desired. The list of occurrences in each category can then be distributed to the appropriate NCC staff member and examined for significance; a record of the global context for each such event can be obtained from the single printed log.

3.7.4 IMP Software Maintenance

The IMP software included the dapability of obtaining a reload of its core memory, except for a small region containing constants specific to each machine and the loader itself, from a neighbor IMP under command from DDT or from console switches; the watchdog timer also used this mechanism. It seemed reasonable to use this mechanism for all new software releases once an IMP was installed at BBN. The procedure of releasing new software could then be to mail a tape of only the constants and loader area of

the IMP to each site and at a pre-scheduled time have each site load this tape into the IMP. The entire new varsion of the IMP program would be read into the BBN IMP through the high-speed paper tape reader. The sites adjecent to BBN would then direct their IMPs to reload from the BBN IMP. Once this was completed, sites adjecent to the reloaded sites would, in turn, reload from these sites until the new software was propagated through the entire network.

Using this software release procedure, site personnel were involved in the release only at those times when the constants and loader area, or the IMP=to=IMP communications protocols, were changed; in many cases new software could be released using only the BBN paper=tape reader and DDT. Thus the software development staff began to place emphasis on keeping the software "compatible" across IMP releases; that is, insuring that the protocol which the IMPs use to communicate with each other remain unchanged.

Finally in mid=1972, when the network had grown to about 30 nodes, it became necessary to release an incompatible version of the IMP software, It proved impossible to coordinate a release based on the paper tape procedure, and therefore it was necessary

to devise a new plan which would not require ĝηγ site intervention. In spite of the incompatible nature of the new software. The plan which evolved required several staps: first, DOT was used to overwrite the logder area of each machine with a new loader which could communicate either with the old version of the softwere or the intended new version (a bit in the IMP=te=IMP header was usurped to differentiate the systems). Once this step had been completed at every IMP, one of the dircults was disconnected from the BBN IMP and connected instead to a development machines as illustrated in Figure 3+7. This second "88N IMP" was running the new IMP software, but this was not confusing to the maghines running the old systems the software incompatibility meant they could communicate with only one "88N IMP" rather than two. A message was then sent from the console Teletype of the "old" 88N IMP to the DDT of site "A" (refer to Figure 3-7) instructing site A to reload (taelf from the "new" 88N IMP. When this reloading process was completed, there was a "natwork" of two machines running the new system, and a network of the remaining machines running the old system. This procedure was repeated for each machine until the entire network except for the "old" BBN IMP was running the new system. This IMP was then disconnected from the network and its circuit(s) connected instead to the "new" ABN IMP,

Figure 3-7: Incompatible Software Release - Topology Example

An additional complication remained to be faced, the case where the topology did not allow two paths from 88N to an IMP, one through the old system and one through the new system. Again referring to Figure 3=7, the "arm" containing two machines (C and D) has only one connection to the rest of the network (through

8). The procedure for this case was as follows. First, the NCC insures that the new additione is operating at IMP E. Next, using DDT, a software "patch" is installed in IMPs B. C. and D. such that each patched machine will ignore relead requests from its neighbors. Then IMP D is instructed to reload (tself from IMP C; similarly, IMP C (a instructed to reload itself from IMP B. Finally, IMP B (s instructed to reland itself from IMP E. Since IMP E (a not patched to ignore reload requests, it immediately transmits a copy of the new software to IMP 8. This reload overwrites the patch which told IMP B to ignore reload requests, so IMP B will then honor the next of IMP Cfa (repeated) reload requests. Similarly, this reload overwrites the patch in IMP C, and thus the new software is propagated to IMP O.

Of course, with the advant of the TIP, even supposing all the TIP programs remained (dentical, it was not possible to quarantee that a Terminal IMP would have any neighbor which was a Terminal IMP (of course, the "IMP" portion of a TIP remained identical to other IMPs). Thus, in the event of any changes or failures in a Terminal IMP, a mechanism to reload the Terminal IMP software, without reliance on a neighbor, was required.

The first approach to the problem of reloading TIPs was to load the TIP software, which might be specialized for each individual TIP, into the "development TIP" (which was used for stand-slone software checkout of TIP software), connect this TIP to the network, and install in it a small routine which could interact with the DDT at the target TIP site. This routine read a word of TIP code and sent this word to the DDT at the remote site with instructions to load it into memory at the same location from which it had been taken. The process was repeated until the entire TIP program had been loaded (through the network) into the remote TIP.

Although this process worked, it was sumbersome and required a great deal or operator activity with the resultant possibility of error, as well as interrupting software development work which might be taking place. Thus, as the number of TIPs in the network increased, pressure mounted for some more automatic method of reloading. It seemed obvious that a reliable Host computer, equipped with a bulk storage unit, could store a copy of the program for each TIP. Even if the TIPs became very divergent, it would not require a great deal of bulk storage to store a separate copy for each TIP. By late 1972 facilities which supported this activity had been implemented on a POP=1

Eventually these support facilities, as well as others which had been implemented on the PDP=1, were moved to TENEX systems because of their greater power and the reliability inherent in their redundancy. These facilities followed the carrier procedure, as described above, of using DDT to load a copy of the TIP program, or a patch to the IMP or TIP software, one word at a time in what amounted to a simulation of an operator's typing. However, this method could not be extended to the loading of IMP software either for failed mechines or for new releases, since use of DDT through the network requires the entire set of IMP software to be running.

As time pessed, and the network became increasingly large and varied, the mechanisms of complete memory reload from a neighbor and of selective reload via "sutemated" DOT became increasingly awkward and inadequete, Several examples of needs not well satisfied by those mechanisms are:

me The NCC wanted to be able to propagate new releases by sending simple commands to each target IMP telling it to accept core transfers from some other designated IMP.

- To examine an IMP that oranhed and get it back on line quickly the NCC needs to be able to ship the core (mage of the failed IMP back to the NCC to be exemined later by a programmer.
- A general reloading schame was desirable to load "add=on" pledes such as the TIP, VDH, or Satellite IMP code, which are not part of the basic IMP system, and Hosts for which the NCC has full responsibility such as the PLI.
- with the addition of Pluribus IMPs to the Network, the existing method of loading an IMP from its neighbor was insufficient. Once a Pluribus IMP has been installed, most network configurations require a general way to load a remote IMP across the network through dissimilar machines.
- Loading or dumping a complete core image of a Satellite

 IMP using a broadcast satellite link is not convenient

 since it prevents use of the link by other nodes for

 several seconds,
- A log(ca) adjunct to sending core (mages through dissimilar machines is sending core (mages through

(perhaps identical) machines which may be running different software. Thus, a new compatible software release could be accomplished in a more modular fashion by loading individual IMPs at the most appropriate time for them, without being constrained to first load at least one neighbor IMP.

The mechanism constructed for accomplishing all of the above is a packet core transfer mechanism which consists of two endpoints == the font (source of the core image) and the sink (destination of the core image) == and a protocol governing the reliable transmission of the packetized core image. The font puts pieces of the core image into packets and sends them over the network; the sink receives packets from the network and builds the appropriate core image.

Either font or sink can take one of three basic forms. The first is that of a real natwork Host, such as a TENEX used by the NCC. Thus the TENEX could be the font for sending "addwon" (i.e. TIP, VDH, etc.) pieces to running IMPs, or the sink for receiving core images of running IMPs for periodic (preventive maintenance) verification of correct software. The second form is that of the "block" fake Host in a running IMPs. This fake

Host is the font for any transmission of core from the running IMP (i,e, for verification at the NCC or for loading another dead IMP of a similar machine type), and the sink for any pieces for the core image loaded into the live IMP (i.e. "add=on" pieces). The third form is a standmalone program running in a nonwfunctioning IMP, in conjunction with a live neighbor IMP's block fake Host, which acts as the port into the network for the failed IMP, and is explained more fully below. The block fake Host is able to distinguish between packets addressed to itself (form two) which require it to load or dump its own IMP's core, and packets addressed to a failed neighbor IMP (form three) which it reformats and sends over a direct circuit from itself to that neighbor.

To bring up a non-functioning IMP, the required stand-elone program (that may have to be loaded on-site) should be as short and simple as possible. Therefore, it should know nothing about network procedures such as acknowledging, routing, RFNMs, etc., and talks over its modems in "blocks", a data format which is different from any other type of pasket, and which cerries much less header information. The live neighbor IMP at the other end of the modem also understands blocks. The block fake Host converts packets striving in standard forms for its

non-functioning neighbor into blocks and sends them over the modem to that neighbor. It also converts blocks arriving from its non-functioning neighbor into packets which are then transmitted over the natwork. When operating in this third form, the block fake Host is completely unaware of the nature and extent of the core transfer; thus, the internal structure of the failed IMP and its neighbor can be completely dissimilar.

The protocol calls for an initiation handshake, followed by one-way flow of pieces of dore image, with possible resynchronization messages in the reverse direction, followed by e termination handshake. The initial request for a core transfer may originate at the NCC, or in a healthy IMP which has a dead neighbor, or from a dead or dying IMP, and goes to the font. The font is now open and ready to send core; and sends a request towards the sink, which then becomes open and ready to receive The font follows the Initial request sequence-numbered pieces of the gare image, and then a "last message" message. If piece-numbers ever get out of sequence, of a piece is excessively tardy, the sink sends back a piece-number reset request, and waits for the correct piece to arrive. When the sink Gets the "last massage" message, it echos it back to the font, causing the procedure to be terminated.

3,7,5 Hardware Maintenance

By m(d+1972 the Natwork Control Center staff had been expanded to provide operator coverage 24 hours a day, seven days a week. At about this time, first order responsibility for network operations was assigned to the operators, with back-up provided by the hardware and agityare dayelopment staff membera who were instantly available at the NCC during the working day and available at home after hours. The operating staff also assumed responsibility for a number of peripheral issues such as the scheduling of installation of new mitem and new circuits, the coordination of routine maintenance activities so as not to disrupt network connectivity, and so forth. Field maintenance personnel from both the common carriers and the Honeywell repair offices (who had primary responsibility for the repair of IMP hardware problems) were instituted to accept maintenance calls only from the Network Control Center, rather than from the individual sites, since only the NCC had a global picture of the network and was able to see how problems in one area might be related to problems in another area.

As a sidelight, even prior to this time the carrier personnel had been instructed to accept trouble calls only from

BBN. As the network grew into new areas, far removed from Cambridge, this frequently led to difficult interactions with carrier personnel, who could not easily believe that a group in Cambridge could be accurately reporting the state of a circuit rumning between, for example, Urbana, Illinois and Sait Lake City, Utah, Fortunately, all circuits were ordered by the government through a gentral office in Washington, and this office was able to observe the NCC track record in reporting failures in other areas of the country, When disputes arose involving a new circuit, a cell from the gentral facility was usually sufficient to convince the repair office that the Network Control Center's repair request was believable. It is important to note that this believability was the direct result of the great core that had been taken to build adequate diagnostic tools into the IMPs, so that the carrier personnel were not called because of problems which were actually caused by the IMP hardware or software.

In early 1974, a hardware engineer with a considerable amount of programming experience was added to the Network Control Center Staff and given responsibility for the maintenance of operational IMPs. The maintenance engineer initiated a quality assurance program which is carried out on each new IMP and each

new interface before it is permitted to be shipped from 88N to the field. The testing procedure is sufficiently routine so that the majority of it can be carried out by the NCC operators during times when they are not occupied with more urgent tasks.

However, a key result of assigning a maintenance engineer to the NCC was the discovery that, with the sid of powerful Host computers, and a rapid communications system, and the use of IMP tools such as DDT which allowed routines to be added to a running IMP to test aspects of its hardware performance, it was possible to maintain IMP hardware by debugging, rather than by trouble shooting. Further, maintenance and debugging is actually of a higher quality, because it is controlled by experts who can be concentrated at the center rather than by the lass highly trained personnel normally available on a field maintenance staff. Of course, the field maintenance staff is still necessary when physical presence is needed at a field site; e.g., when the central personnel wish a herdware component changed.

One example of the changed approach to maintenance is the study of memory errors. The trouble-shooting approach to this problem is to swap memory driver cards from one memory stack to another to see which, if any, causes the problem to move. This

typically requirem a lot of onesite time and many service interruptions. The debugging approach uses the packet core mechanism plus a verification program running on TENEX transfer which retrieves a dopy of the core memory of a running IMP and compares this, word by word, to an image of correct memory contents on the TENEX mass storage. Using this method it is possible to identify any pattern of dropped or picked bits and from this pattern the particular driver in error. example of a diagnostic procedure is the use of DDT to accurately measure an IMP's resistine clock. Failure of the real-time clock can lead to several problems with an IMP, particularly to too frequent or too infrequent probes of the circuits to which the IMP is connected, which is likely to result in one of the IMPs connected to the circuit declaring it unusable. The test of the clock can identify it as the cause of this trouble much more easily than can trouble-shooting in the field.

Use of these centralized techniques leads to a different economic balance between the central staff and the field staff. Thus, by 1975 it became desirable to discontinue Honeywell hardware maintenance responsibility and add this task to the NCC in order to "trade" field staffing support for central staffing support. The result of these changes in maintenance philosophy

was a reduction in the average hade downtime due to hardware and softwere fallures from about 1.5% to about 0.25%, a reduction of about a factor of six,

REFERENCES

- Abremson70 = N. Abremson, "The ALOHA System Another Alternative for Computer Communication", AFIPS 1970 National Computer Conference Proceedings, Vol. 37, November 1970, pp. 281=285.
- AFCS72 Air Force Communications Service, "AFCS Test of the ARPA Network", Richards-Gebaur AFB, Missouri, January 1972.
- Bouknight73 = W.J. Bouknight, G.R. Grossman, and D.M. Grothe, "The ARPA Network Terminal System = A New Approach to Network Access", Proceedings of the Third Data Communications Symposium, November 1973, pp. 73=79.
- Braden73 = R.T. Braden, "Update on NETRJ9", Network Working Group Request for Comments #599, December 1973.
- Butterfield74 = S.C. Butterfield, R.D. Rettberg, and D.C. Walden, "The Satellite IMP for the ARPA Network," Proceedings of the Seventh Annual Hawaii International Conference on System Sciences, Honolulu, January 1974, pp. 70-73.
- Carrio = S. Carr, S. Crocker, and V. Cert, "Hostwhost Communication Protocol in the ARPA Network", AFIPS 1970 National Computer Conference Proceedings, Vol. 36, May 1970, pp. 589-597.
- Cerf74 = V. Cerf, and R. Kahn, "A Protocol for Packet Natwork Interconnection," IEEE Transactions on Communications, Vol. COM=22, No. 5, May 1974.
- Cerf76 * V. Carf, A. McKenzie, R. Scantlebury, and H. Zimmermann, "Proposal for an International End to End Protocol," ACM Computer Communication Review, Vol. 6, No. 1, January 1976, pp. 63-89.
- Cosell75 * B. Cosell, P.R. Johnson, J.H. Malman, R.E. Schantz, J. Sussman, R.H. Thomas, and D.C. Walden, "An Operational System for Computer Resource Sharing", Proceedings of the Fifth Symposium on Operating System Principles, Austin, Taxes, November 1975, pp. 75-81.

- Dorin77 = R.H. Dorin and D.E. Eastlake, "Use of the Datacomputer in the Vela Saismalogical Network", Proceedings of the IEEE International Symposium on Computer Aided Saismic Analysis and Discrimination, Falmouth, Masachusetta, June 1977, pp. 35+39.
- Englebart73 = D,C, Englebart, R,W, Watson, J,C, Norton, "The Augmented Knowledge Workshop", AFIPS 1973 National Computer Conference Proceedings, Vol. 42, June 1973, pp. 9=21,
- Falk76 = H. Falk, "Reaching for a Gigaflop", IEEE Spectrum, Vol. 13, No. 10, October 1976, pp. 65=69.
- Forgle75 of J.W., Forgle, "Speech Transmission in Packet-switched Store-and-Forward Networks", AFIPS 1975 National Computer Conference Proceedings, Vol. 44, pp. 1370142.
- Forsdick?? = H.C. Forsdick, R.E. Schantz, R.H. Thomas, "Operational Systems for Computer Networks", BBN Report No. 3614, July 1977.
- Henderson77 = D.A. Henderson and T.H. Myer, "Issues in Message Technology", Proceedings of the Fifth Data Communications Symposium, Snowbird, Utah, September, 1977.
- Kahn75 ~ R.E. Kahn, "The Organization of Computer Resources into a Packet Radio Network", AFIPS 1975 National Computer Conference Proceedings, Vol. 44, pp. 177=186,
- Kehi73 = W. Kehi, "The UCLA Campus Computing Natworks An ARPANET Resource", Educom Bulletin, Vol. 8, No. 4, Winter 1973, pp. 10-17.
- Maril175 = T. Marili and D. Stern, "The Datacomputer = A Network Data Utility", AFIPS 1975 National Computer Conference Proceedings, Vol. 44, pp. 389=395.
- McKenzie72 w A.A. McKenzie, B.P. Cosell, J.M. McQuillen and M.J. Thrope, "The Network Control Center for the ARPA Network", Proceedings of the 1972 International Conference on Computer Communication, Washington, D.C., May 1972, pp. 189-191.

- MoKenzie75 w A.A. McKenzie, "The ARPA Network Control Center", Proceedings of the Fourth Deta Communications Symposium, Detaber 1975.
- Matcalf#71 = R. Matcalfe, "Some Historic Moments in Networking", Network Working Group Request for Comments #89, January 1971.
- Mobley77 = R₄L_a Mobley, "Comments on the ARPANET", personal correspondence with Bolt Beranek and Newman, August 1977_a
- Ornstein72 w S.M. Ornstein, F.E. Heart, W.R. Crowther, H. X. Rising, S.B. Russell, and A. Michel, "The Terminal IMP for the ARPA Computer Network", AFIPS 1972 National Computer Conference Proceedings, Vol. 40, pp. 243-254,
- Panko77 = R.R. Panko, "The Outlook for Computer Mail", Telecommunications Policy, June 1977, pp. 242=253.
- Retz75 * D.L. Retz, HELF * A System for Network Access*, Proceedings of the IEEE Intercon Conference, New York City, April 1975.
- Retz76 = D.L. Retz and B. W. Schafer, "Structure of the ELF Operating System", AFIPS 1976 National Computer Conference Proceedings, Vol. 45, pp. 1007=1015.
- Robinson76 R.A. Robinson, "The National Software Works: Its purpose, Status and Future Plans", Rome Air Development Center, briefing to ARPANET Spansors Group Meeting, November 1976.
- Thomas 73 * R.H. Thomas, "A Resource Sharing Executive for the ARPANET, AFIPS 1973 National Computer Conference, Vol. 42, pp. 155*163.

4. OBSERVATIONS

For many of the people in government, at the major contractors, and in the participating universities and research centers the development of the ARPANET has been an exciting time which will rank as a high point in their professional careers, In 1969 the ARPANET project represented a high risk, potentially high impact research effort. The existence of the net in practical useful form has not only provided communications technology to meet many short term needs, but it represents a formidable communications technology and experience base on which the Defense Department as well as the entire public and private sectors will depend upon for advanced communications needs. The strong and diverse experience base generated by the ARPANET project has placed this country shead of all others in advanced digital communications science and technology.

4.1 Social Issues

Somewhat expectedly, the network has facilitated a social change in the United States computer research community, it has become More convenient for geographically separated groups to perform collaborative research and development. The ability to easily send files of text between geographically remote groups, the ability to communicate via messages quickly and easily, and the ability to collaboratively use powerful aditing and document production facilities has changed significantly the "feel" of collaborative research with remote groups. Just as other major improvements in human communication in the past have resulted in a change in the rate of progress, this social effect of the ARPANET may finally be the largest single impact of the ARPANET development.

A non-trivial question of sonsiderable importance to the country is why has the ARPANET project been so successful. It would certainly be nice if the same formula could be applied to other pressing national needs. While timing, accident, and luck must not be discounted, it is possible to identify several possible contributing factors:

o Within the Defense Department, and even within the government as a whole, ARPA has had an unusually high degree of freedom in identifying the right people and then funding those people in a relatively rapid and

efficient manner. Further, even within ARPA, Larry Roberts was an unusually gifted individual with a special blend of Personal technical competence, political persuasiveness, and drive. Although there are always some constraints, Roberts was able to manage the research and development at the many organizations involved with the network in a very unfettered fashion. It was clearly benevolent dictatorship rather than committee management, although Roberts "listaned well" before deciding on a course of action. Strong management also existed within the organization of the prime contractor, Soit Beranek and Newman Inc., where project control was vested in a single strong individual. It is probably the case that only ARPA could have built the ARPANET with the same degree of SUCCOSS,

Despite the fact that the ARPANET was a government project and further despite the fact that it was run within the Defense Department, it was possible, at least initially, to free the research and development from some of the constraints which often seriously hamper other activities. In particular, the project was entirely unclassified and even the interactions with some highly classified agencies managed to stay on an unclassified beets. Second, the provision of network

service was for a very long time provided as a "free good"; this allowed people to experiment with the use of without being forced to make early and the metwork probably (1) advised gost/benefit decisions. despite a deeply ingreined government and Department worry about unauthorized use of dovernment facilities, it was possible to build the ARPANET without complex administrative control over access or complex lagin procedures or complex accounting of exactly who was using the net for what. Finally, the ARPANET did not have to interconnect with other existing end/or decrept communication systems; it was possible establish line protocols and interface standards nove, in the best ways that could be devised. Thus, the ARPANET program was incredibly free of "artificial" requirements and was able to consentrate intensively on primary required research and dayelopment. Much later in the project, after success had been assured, it was relatively easy to reopen consideration of some of these issues; and at the current time, for example, the Defense Communications Agency is charging user groups for network access, and there have been experiments (using private line interfaces) in the transmission of classified data over the ARPANET, and research was done on how to implement network login procedures, etc., etc. The key was to evoid such extra confusions at the outset of the project,

A very convenient fact was the common ARPA support of both the "network authority" and the initial early group of network users. It was possible for ARPA to strongly encourage a cooperative attitude and cooperative engineering at the time in the project when such cooperation was most critically necessary.

In sum, the project was an illustration of what can be accomplished with strong technically sophisticated central management, adequate resources, and a clear headed undeviating concentration on the central research and development issues.

largest single surprise of the ARPANET program has been the incredible popularity and success of network mail. There is little doubt but what the techniques of network meil developed in connection with the ARPANET program are going to sweep the country and drastically change the techniques used. for intercommunication in the public and private sectors, **5** y hindsight, one can easily see the reasons for this success. primary prior existing communications techniques (the U.S. poste) service and the telephone) have extremely serious deficiencies: The postal service has become more expensive and (ts. performance has dropped; one must expect delays measured in days for letters, and some times the delay is weeks. In the case of the telephone, our increasingly mobile society makes it difficult to reach the desired person, and reaching 10 distributed people within a short per(od of time is essentially impossible. Leaving telephone

works if a careful system is established and **меввасе** а travelling individual sheeks beak with his answering service or secretary frequently, but it is often inconvenient and is usually limited to the prime shift working day of the secretarial world. To find a busy person able to edgept a phone call at the time you make it is truely unusual, and some officials have desks covered with little pink slips reporting on telephone messages with requests for return calls. Into this milieu was dropped a technique of natwork mail where at any time of the day or night one can send a message to any number of other people and expect that message to semiminatenteneously be available in the computer mailbox of all the recipients. Then one only has to essume the habit on the part of all individuals using the system to occasionally check their mailboxes when they are free and not at and the performance of the dommunications meeting. 1 8 immeasurably improved over the postal service or the telephone. With the addition of sophisticated tools for answering messages, filing measages, forwarding measages, and categorizing measages, the system takes on yet another step function of performance over the alternate possibilities. In the space of lust a couple of years this "computer center curiosity" became a smashing success on the ARPANET; there was a sufficiently large community of individuala who wished to communicate, and they all were relatively mobile, and the system overnight become a way of life, The implications of this kind of success are enormous. Perhaps advances would have eventually come anyway as the U.S. such

Postal Service cost-performance continues to drop, but there is no question that the ARPANET program provided a truely convincing demonstration of the power of this approach. The change will not be overnight, because it does depend upon the availability of terminals accessible to wider and wider populations, but commercial systems are already available and substantial DoD experiments are already under way.

Before leaving the class of rather broad observetions, it is perhaps appropriate to note one disappointing political phenomenon that was visible in the ARPANET program. In partigular, the ARPANET program provided another example of the difficulty of strong research cross=fertilization between various major branches of the federal government. As the success of the ARPANET become clear, many individuals realized that there were other communities of research in the United States who could substantially gain from some experience with the ARPANET. For example, the research communities supported by the National Institutes of Health and the research communities supported by the National Science Foundation could have gained substantially by a serious interaction with the ARPANET program. Based on this obvious feet, several effempts were made to create such cross ties (and, in fact, one important NIH=supported feeility was tied to the ARPANET and did provide some benefit in this direction), However, it was extremely difficult to make progress down this path. A clear "not invented here" phenomenon existed in some quarters, a feeling that DaD programs should not be expanded to include other agencies existed in some parts of the DaD, and in some other agencies a feeling existed that "military" programs should not be mixed with "clean" scientific programs. Further, the uncertainty over the direction in which technology transfer might take place from ARPA did not make it easy for other agencies to commit funding for integration attempts without any clear expectation about the duration of the ARPANET life. A philosophical view of this difficulty might hold that only about a half decade delay in the use of the technology by other research communities was introduced, and why should thet matter? However, at the time it seems hard to condone such difficulty of technical cross-fertilization. Perhaps there might be some way by which the federal government could annourage greater research cross-fertilization between egencies in selected areas,

DRAFT Observations

4,2 Technical Lessons

Leaving the broad social plane, the ARPANET program provided several technical lessons which are worthy of general comments

III-857

4.2.1 Terminal Handling

Rather early in the ARPANET program, it became clear that terminel access to the net via the main. Host computers was an inadequate approach; many classes of users required direct terminal access to the net in order to use major Hosts at other locations. The first response to this pressure was the design of the Terminal Interface Message Processor (TIP) by the ARPANET prime contractor, Bolt Beranek and Newmen Inc. The TIP was destaned t o address a limited problem. to handle character-oriented asynchronous terminals only and to be integral part of the network authority and not available for user programming or special user features. This limited goal and hard-nosed attitude permitted the rather rapid completion of the TIP design, the fielding of many TIPs, and the rapid availability of wide spread terminal access to the network. Thus, the TIP effort was extremely successful in reaching its limited goal.

Unfortunately, but perhaps not surprisingly, the limited goal and absolute restriction on user programming created considerable unhappiness in portions of the potential user community, and created considerable pressure for other "better" terminal access techniques. Some of the complaining was fully justified; this approach to servicing interactive character terminals had "leapfragged" a whole segment of the industry that was congerned with batch processing and the use of synchronous line disciplines from such batch processing units. Other

complaints were less rational and more self serving, where some groups really wanted "a computer of their own" under the guiae of obtaining terminal access to the ARPANET. Still other criticism was based on an honest difference of opinion as to the relative ease of designing and deploying a terminal access device that would permit user programming and would handle a broader class of terminals. Finally, some criticism was provided by highly sophisticated users who were used to the terminal support "services" provided by the most advanced lenga Hosts, and did not like the limited services provided by the tiny minimal resident in the TIP.

In response to this pressure, ARPA for a time supported the development of a device salled the "ARPA Network Terminal System" (ANTS). This was a system based on a PDP=11 which was intended to provide user programming ability, the handling of synchronous terminals as well as asymphronous terminals, and a more powerful set of services to the terminal user. Unfortunately, the goals were somewhat ambitious and, although a few ANTS were put into field, the configuration management of the program, the the difficulty in debugging fielded versions of the SVSt 6m4 delays in implementation eventually led to a cessation of ARPA support for the effort. Since the pressures did not diminish, a second similar attempt was made with a system called "ELF". Here ARPA support was provided to try to standardize, improve, and deploy a PDP=11 based terminal support softwere

system which had been independently developed for a particular Again, the hope was to permit user programming and the handling of a wider class of terminals. From the viewpoint of deployment, the efforts to use ELF were much more successful; goals were far more limited, and more attention was paid to orderly development and maintenance. However, ELF has probably been more useful as a base for individualized miniwhosts, serving local specialized Host functions than as a means of widely distributed network access for large numbers of terminal users. the time the ARPANET was transferred to DCA, primary terminal access was still either through the main metwork Hosts or via the many TIPs in the network. There are several morels to this history. First, it is extremely difficult to build a system which can handle all possible terminals; it is a bit like the "everything" machine and leads to an unlimited expansion of the problem. Second, very great attention must be paid to configuration management, central program design and release, central maintenance and software support, and close control of program changes if an evolving computer-based device is to be deployed in considerable numbers around the network. There 18 still an open technical question whether such a device could be sensibly and costweffect(vely fielded in a way which would permit both some level of user programming to tailor the functions of the device and a standard set of controlled basic services; at this point in the development of the technology a conservative approach precludes such user programming if a device is to be widely deployed,

Another related aspect of the terminal handling problem has to do with the management extent of the network authority. It was discovered that when an unsophisticated user typed at terminal in a remote location and something went wrong (that is, the proper response was not rede(ved), the user typically "blamed the network" despite the fact that any number of possible things: the terminal itself, the local modem, the local line, the modem at the other end of the local line, the terminal handling device (e.g., TIP), the network (twelf, the eventual Host computer, or any similar point on the return path; could have been at fault, Thus, it is extremely important that the network authority have adequate administrative control over the entire collection of equipment from the terminal right through to the Host computers it is to be in a position to respond to such unstructured outcries of rage from the end terminal user. Ouring the ARPANET numerous "crisis" incidents were presipitated by the Orograma failure of equipment which was not in any way under the control the network authority. [Perhaps the best single example was difficulties in the local on-base twisphone system at Ames took a massive effort to eventually uncover the offending equipment and where the network authority was forced in (is own defense to participate and lead in the massive debugging activity even though the offending device eventually found was clearly outside the network authority's control. The generalizable lesson from this kind of history is that devices that attach to networks must be extremely clearly in <u>somebodyfs</u> camp

responsibility. Either the device must glearly belong to the network authority and must have builtein techniques for debugging and failure location, or the device must clearly belong to the Host organization or some other well (dentified group who understands its responsibilities for maintenance and trouble location; devices cannot sit in cracks between different authorities. Although this idea is really quite simple, it is frequently overlooked.

A final point on the general terminal handling problem has to do with the location of the Mintelligence for dealing with terminal related matters. In general, the data processing power can either be in the terminal itself, in the terminal handling network port (e.g., the TIP), in the eventual final Host computer or, of course, distributed among these three locations in some way. As the cost of data processing is dropping, more and more intelligence is appearing in the terminal itself and this entire matter must be periodically reviewed to ensure that the system performance can take advantage of technological change.

4,2,2 Reliability and Fault Isolation

At the outset of the ARPANET design considerable effort wear expended toward builtwin techniques of fault isolation and recovery. For example, IMP=to=IMP circuits could be cross patched for fault isolation, IMPs themselves had power fail recovery mechanisms, and mechanisms for reloading one IMP from a neighbor. As the network grew, it was therefore somewhat of a surprise that the initial efforts down this path were not nearly enough. Over the life of the ARPANET program, there was a nearly continuous effort to add techniques and improve existing techniques for fault isolation, rapid recovery, and containment of failures. The lesson probably is that it is difficult to have too much attention to fault isolation and recovery mechanisms. Many interesting techniques were aloaly added to improve the network's ability to cope with troubles.

for critical places of code such as the routing computation, the code (tself was checksummed before it was operated. Similarly, key data structures were checksummed before they were accessed. This kind of protection was added when it was discovered that the trouble in these particular classes of computations could cause network-wide feilures rather than simply local difficulties.

- It was discovered that dislein modems for remote terminals could break or "hang" and there was no simple way to discover this had happened. A technique was designed wherein a centrally located autodisler controlled by the network authority used an out WATT's line to periodically dislend test all dislein ports all over the network (at least in continental United States).
- It was discovered that when a difficulty occurred in an IMP which caused an automatic program reload, the necessary debugging information was lost and the trouble would likely redure. A technique was added to automatically dump offending code before a reload was attempted which then permitted comparison with a master copy and improved cability for debugging.
- The ARPANET program was forced to cope with the simultaneous need for a twenty-four hour a day continuously operating reliable system and at the same time relatively constant levels of growth, modification, and change. After some early false starts when large changes introduced periods of substandard performance, a technique was evolved whereby greater attention was paid to dividing changes into small incremental changes that were compatible with the previous system. Then, when

small incremental change or retreat taken to the previous release.

As network users began to depend upon a few particular "service" Hosts for a wide variety of services, including message services, it became more important that such Hosts maintain availability to the network, In some cases, a Host might be operating adequately as perceived by its own local operators and yet in some way not be properly servicing the network connection. A technique was evolved whereby the network authority added software tools to "watch" these particular critical Hosts and was then in a position to urge corrective action by the local Host authorities if and when necessary.

4.2.3 Maintenance Management

In the early years of the ARPANET program, the IMPs and TIPs of the network were maintained by aubsontract to the manufacturer of the basic mini computer (Honeywell). This represented a costmeffective domorage baceuse Honeywell had maintenance facilities in many cities. However, this rather standard form of computer maintenance was aimply inadequate for the reliability requirements the ARPANET, After a o f time, maintenance of the network nodes was undertaken directly by Solt Berenek and Newman Inc., the network prime contractor, and special new techniques were developed which greatly improved the average nodel reliability. In particular, techniques were developed of central maintenance management wherein a very strong teem of hardware and software experts at the central Network Control Center acted in close condert with a small number of field personnel who became highly expert in the IMP and IIP machines. The contral staff could use the natwork (tself to observe the behavior of an offending node and it could talk difficulties with the local maintenance angineer. Further, the people in the field became much more dedicated and responsive to ARPANET difficulties as compered to time-shared Honeywell personnel who had to take responsibility for many different kinds of equipment in their geographical territory, This approach to maintanance probably has important benefits for other distributed systems, especially as those I I w emersys

Section 4

DRAFT Observations

increasingly be interconnected in networks and thereby accessible to central maintenance,

4.3 Impact on Other Research Areas

A final class of observation should appropriately deal with the extent to which the ARPANET program made it possible for other lines of research to be affectively pursued. There have been at least three different (moortant ways in which the net has had this kind of impact: (1) The encouragement of research on large specialized resources which would not have been cost-effective without a network as a distribution mechanism; (2) the encouragement of computer research on the aublect of interprocess communication which would be generally useful even fin. without a network, but becomes critical the (3) the development of interpersonal environment: and communication techniques which again are generally useful without a network, but which become critical in a network environment.

4,3,1 Spec(alized Resources

Several important research projects with considerable potential for the Defense Department would simply not have been undertaken without the network's existence as a distribution tool. One example is the "National Software Works" where ettempts are being made to develop integrated sets of computer tools for program design and construction which then can be used by remote groups which do not have such construction tools locally available. A second example is the Datacomputer project where a large store and efficient techniques for accessing this store were developed with the notion that this store would be used by diverse groups remotely distant from the store itself; In each dase, the development of such techniques may be generally useful, but the development would not have been undertaken if the initial utilization was forced to be an a local basis.

4.3.2 Interprocess Communication

Very early in the ARPANET project, one remearcher made a comment which can be parephresed as why bother trying to get programs to intercommunicate over a network when we don't even know how to make programs intercommunicate within a single Host". The answer (a that there was not sufficient intentive to work on interprocess communication in a deep way until the existence of the network made such interprocess communication more critically necessary. Here is a case where the network has forced research in areas which will be useful even aside from network applications. Both at a formal and a practical level considerable research has now been undertaken on interprocess communication and some yaeful results have been obtained.

4,3,5 Person-to-Person Interactions

In the case of the spectacular success with network mail but also in some other cases of somewhat lesser impact (such as document production facilities, editing facilities, etc.), the existence of the network dreated a much larger community of individuals who had need to intercommunicate and thus considerable effort has been expended on such interpersonal communication tools. Such work will be of benefit in situations even without networks and will certainly be of benefit as networks become more a way of life.

Section 4

DRAFT Observations

4.3.4 CONCLUSIONS

These kinds of benefits have represented a very important contribution of the ARPANET program and it is somewhat fitting to end on the note that the ARPANET program has had a strong and direct feedback into the support and strength of computer science, from which the network itself sprung.