# A PROTOCOL FOR HIGH SPEED OPTICAL NETWORKS
# FROM A BASIS OF SATELLITE PROTOCOL DESIGNS

by

## ROGER KIRK ALEXANDER

B. Eng., Electronic and Electrical Engineering

University of London (1987)

SUBMITTED TO THE DEPARTMENT OF

ELECTRICAL ENGINEERING AND COMPUTER SCIENCE

IN PARTIAL FULFILLMENT OF THE REQUIREMENTS

FOR THE DEGREES OF

MASTER OF SCIENCE IN ELECTRICAL ENGINEERING

and

MASTER OF SCIENCE IN TECHNOLOGY AND POLICY

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

June, 1991

Signature of Author.........................................................................................................
Department of Electrical Engineering and Computer Science
May, 1991

Certified by.........................................................................................................
Pierre A. Humblet
Thesis Supervisor

Accepted by.........................................................................................................
Richard de Neufville, Chairman
Technology and Policy Program

Accepted by.........................................................................................................
Arthur C. Smith, Chairman
Committee on Graduate Students
Department of Electrical Engineering and Computer Science

# A PROTOCOL FOR HIGH SPEED OPTICAL NETWORKS
## FROM A BASIS OF SATELLITE PROTOCOL DESIGNS

by

## ROGER KIRK ALEXANDER

Submitted to the Department of Electrical Engineering and Computer Science
on May 10, 1991 in partial fulfillment of the requirements for the Degrees of
Master of Science in Electrical Engineering and Master of Science in Technology and Policy

## ABSTRACT

The ability to utilize the vast transmission bandwidth of optical fibers rests in designing communication protocols that can facilitate a very high degree of information multiplexing. With very large transmission capacities andspeeds a major characteristic of these optical fiber networks becomes the very high ratio of end-to-end propagation delay to message transmission time. This, together with the multichannel architecture that results from the need to overcome the speed limitations of electronic communication, presents unique challenges for protocol design.

The significant feature of satellite communication protocols lies in their ability to operate in an environment where the ratio of signal propagation time to transmission time is very large. Other characteristics such as the broadcast, multiaccess nature of satellite communication provide a number of technical similarities with the situation encountered in high speed optical network systems. These operating characteristics provide a basis for establishing potential approaches to the problem of protocol design for optical networks.

The first segment of this study was thus to obtain a clear classification and categorizing of satellite protocol methods, and an understanding of their performance capabilities under various operating conditions. This was achieved through an exhaustive survey of the many schemes that have been proposed. What resulted was the establishment of the common bases from which a number of the schemes are derived as well as an understanding of the elements that comprise the most successful approaches.

Together with an examination of the work presently done in the area optical network protocols, the understanding gained from the study of satellite protocol designs was used as a basis for initiating studies aimed at achieving an effective protocol scheme for optical network communication. Given a number of practical constraints imposed by the multichannel optical architecture, an adaptive, random access based scheme was designed. This protocol, the Extended URN (XURN) scheme, is based on the ideas of the URN scheme introduced by Kleinrock and Yemini [KY 78] for satellite operation. The method provides simplicity of implementation and addresses a number of the short coming of some of the early optical schemes as well as issues of practical design. The performance is also superior to other optical schemes presently proposed.

In the final segment of this thesis we take a step backward and look objectively at the overall process involved in this study to better understand some of the internal factors that influence technology innovation in areas of focussed technical research. We draw extensively on the wide literature of innovation process research to highlight the course that our technological design has taken and how this may have influenced the results obtained. This analysis allows us to make future technology policy recommendations for similar undertakings.

Thesis Supervisor:  Dr. Pierre A. Humblet
      Title:  Associate Professor of Electrical Engineering

## ACKNOWLEDGEMENTS

# Contents

# Chapter 3 - Survey of Proposed Satellite Protocols

# Chapter 4 - Examination of Multiaccess Optical Network Protocols

# Chapter 5 - Proposal for a New Optical Network Protocol

# Chapter 6 - Understanding the Process of Technological Innovation

# List of Figures

# Chapter 1

# Introduction and Overview

## 1.1 Protocol Development for Optical Communication

Fiber optics represents the essential component of future data networks. The optical fibers used in these networks afford transmission bandwidths on the order of tens of terahertz. Since the devices connected across optical networks are presently limited by the speeds of their electronics to tens of gigahertz, is has been realized that the means of achieving efficient utilization of the communication resource lies in the appropriate sharing of the medium. With recent advances in single-frequency lasers and narrowband filter technology, there has been phenomenal progress in the area of Wavelength Division Multiplexing (WDM) technology, as a means of partitioning the optical channel and providing simultaneous, interleaved communication, amongst significant numbers of stations. This approach utilizes the low-loss optical fiber window of about 100 nanometers (on the order of 12500 GHz) to create a high density of multiple communication channels.

The main issue of concern in any approach to the efficient allocation of available transmission bandwidth lies critically in controlling the access to the common communication channel. This must be achieved in a manner responsive to the needs of the individual contending users. In the area of optical network communication, this allocation needs also to be consistent with the requirement of high connectivity as envisaged for future networks which will serve large communities. We can define an

access communication protocol as an algorithm to determine channel access rights, with its objective being to maximize the rate of successful information transmission, subject to some constraint on delay. The implementation of channel access schemes that allow multiple users to share a given resource, have traditionally being driven by the desire for improved performance and reduced unit cost of resource usage. These motivations of better resource utilization through multiple access also apply in optical networks.

With the large bandwidths available, the optical communication environment is no longer one of multiple access to a single channel but rather multiple access to a number of parallel channels interchangeable amongst users. Appropriate coordination and access control in the use of the medium is thus more precisely demanded. If the available bandwidth of these optical networks are to be efficiently harnessed, the emphasis for optical protocols lies not only in appropriate multiple access to individual channels but also in the coordination and control of these parallel, cross linking, multiaccess channels.

The design of such protocols for the control of high speed optical communication must also be framed within the limitations of the architecture and the physical components of a realizable network. As envisaged by Kennedy and Humblet [KH], optical networks of the future might comprise an all optical transparent (broadcast) core within which the flow of signals is not impeded by any optical to electronic conversions. Individual users and other networks will be able to access this core through gateways which provide the interface between the optical network and the electronics of the connected nodes. Given the terahertz optical bandwidths, the all-optical core, supporting thousands of gateways with capacities less than a gigabit per second each, will still not require switching as a means of localizing information. Thus broadcasting will be employed in many applications, as in local area networks.

Due to signal power concerns, however, switching to localize information within the core network may be necessary to limit the number of gateways directly connected to each other. Switching would limit the power reduction experienced by signals put on the network - the power being divided among all directly connected nodes - and so allow more gateways to be supported in the network. The multiple optical channels spaced apart in wavelength would facilitate the simultaneous communication between pairs of stations and across network interfaces.

On each optical link therefore, tunable laser transmitters [CG 88-2] will be used in conjunction with fixed optical filter receivers, or equivalently, fixed transmitters used with tunable optical receivers. While the choice of either of these alternatives does not significantly influence the design of the protocols, the current state of the art in technology appears to favor the fixed transmitters-tunable receivers variant [CG 89]. While a tradeoff exists between the number of frequency channels supported by given types of filters, and their speeds of operation [CG 89], this limitation in device characteristics will not be addressed in the presented protocol concept.

Unlike electronic systems, optical systems do not as yet have the equivalent of high impedance receivers. As a result every coupler introduces a significant power loss. This fact has thus, at present, determined the preferability of a star topology for the network [CDR 90]. The use of a bus or ring topology is less attractive since tandem taps introduce an exponential power attenuation that severely limits the maximum number of taps. The advent of optical amplifiers (such as the demonstrated erbium-doped optical amplifier) may, however, change this trade off. Also, a star topology is more flexible to user growth within the capacity limit of the coupler. The passive star can also be used as a broadcast medium to cross-connect network users. It is this type of topology, therefore, that we shall consider in the designing of a new protocol scheme.

In the star coupler configuration, the hub is a passive transmissive element which distributes information to the periphery. This factor of topology has a direct bearing on the design of the protocols

as it introduces a central component in the path between two communicating stations. Such a central hub in the network, while not precluding distributed control, can in certain protocol applications, contribute to a reduction in packet latency. With the presence of a central hub, the time between packet retransmissions can be shortened to the round trip propagation time between a transmitter and the hub. This is instead of the longer propagation time between a transmitter and a network-connected station receiver [CDR 90]. Also, collisions at station destinations can be avoided by the avoidance of collisions at the hub.

One of the most important characteristics relevant to the design of optical protocols is the long propagation delay across the network. In high speed networks, this delay is long relative to the transmission time of packets, and is becoming more so as the diameter of the networks increase[1]. This long propagation delay has a major impact on the bandwidth allocation techniques and on error and flow control protocols. Carrier-sensing protocols capable of producing throughput levels close to unity are highly inefficient in such large propagation time environments.

In the high speed optical networks of the kind to be considered, each frequency communication channel is representative of the type of transmission medium encountered in a multiaccess, broadcast network as typified by the ALOHA network[2]. The issue of multiaccess control in communicating across the individual frequency channels of an optical is thus one that has already been studied in the context of satellite communication.

In addressing the design of protocols for these optical networks therefore, the examination of the protocols and allocation methodologies employed in satellite communication systems offers an

---

1     With a network diameter of 50 km, the propagation delay is about 250 $\mu$s. With a data rate of 200 Mb/s
      and packets of length of 1000 bytes, the packet transmission time is 40 $\mu$s.
2     The ALOHA network [Abr 70] was developed at the University of Hawaii to provide radio communications
      between a central computer and various data terminals. In the ALOHA network, a single radio
      communication channel exists that is shared, using random access, among all users. This is discussed
      further in Chapters 2 & 3.

appropriate point for initial study. In addition to the broadcast nature of the communication, satellite protocols have in varying forms focussed on the demands of multiaccess communication within the constraints of high propagation delay. Also, as in the case of multi-beam satellites the two-dimensional problem of frequency allocation and channel access has been studied. The comparison with satellite systems is also valid considering that the star coupler and the satellite each act as a central resource.

## 1.2 Introduction to Multiaccess Communication

While a number of multiaccess techniques exist for the establishment of transmission control, the choice of protocol depends upon the environment being considered and the system characteristics to be satisfied. Multiaccess schemes can generally be divided according to three main classes: i) fixed assignment, ii) variable assignment, and iii) demand assignment approaches. These refer to the manner in which the central communication resource is allocated amongst the community of contending users. While any of these approaches could be employed in either a centrally or distributedly controlled manner, the significant propagation delay time encountered in satellite or optical networks generally precludes the use of centrally controlled assignment techniques. In distributed schemes, however, there may be the additional requirement of queue information synchronization where distributed global queues are employed.

In an environment of excess transmission bandwidth, it would appear that the appropriate approach to communication control would be the fixed assignment of individual subchannels to separate pairs of communicating stations or users. While this may be suitable in a network serving a very limited number of stations, the reality of large optical networks with the number of connected stations being far in excess of the number of subchannels, encourages a more dynamic approach to the problem. In addition, computer and terminal data streams are generally *bursty*. That is, over a period of time, the peak data

rate is much larger than the average data rate. This bursty nature of the predominant computer and data communications across the network also means that the use of dedicated channels through static allocation or fixed assignment schemes provide less efficient bandwidth utilization than through dynamically shared approaches. The network and the type of traffic envisaged therefore dictates the use of statistical multiplexing techniques.

With the objective of efficiency in the sharing and utilization of the communication channel, packet communication is mainly considered since it allows available resources to be allocated to users on the basis of shorter time periods. This is achieved by dividing data into small transmission units - packets. However, while packet communication has developed as the natural means of more efficiently sharing the common channel, the operation of the network in a circuit mode may also be of useful consideration. This could allow for the adaptability of the optical network under differing traffic demand scenarios where, for example, voice circuit traffic capability may be required. The ability to accommodate different traffic types would improve the overall utilization of the network.

The ultimate test of appropriate access mechanisms, however, will be the resulting system performance achieved. While multiaccess protocols are evaluated according to various criteria, the performance characteristics that are most desirable are high bandwidth utilization and low message delays. Other attributes are however also important. The ability for an access protocol to simultaneously support traffic of different types, different priorities, variable message lengths, and differing delay constraints is essential, as higher bandwidth utilization is achieved by the multiplexing of all traffic types. The ideal control of the network would be one in which transmission is dynamically adaptable to provide satisfactory performance under varying traffic and throughput demands. In such a model, as traffic conditions varied, the schemes most efficient to the particular conditions could be employed. There would then be a constant adjustment of method to suit to the given traffic demands. The objective would be an access control strategy that dynamically employed the protocol scheme appropriate to traffic levels and statistics.

15

## 1.3 The Central Problem

The problem of multiple access to a broadcast channel is essentially one of decision; which user should have an access right to the channel at any given time (or in the frequency domain; which channel is to be accessed by which user). The decision is also distributed in that each decision maker (user station) possesses only partial state information, knowing only its own needs, insufficient to estimate the overall state. Additionally, the different decision makers need to coordinate their policies. While it is possible to centralize the system's information ([NG 77]) and reduce the problem to one of classical decision making, the overhead and additionally introduced delay in such an exercise generally makes that approach inefficient in distributed data networks such as the ones being examined. Moreover, an algorithm that depends on detailed state information and coordination may be unreliable. The protocol methods employed for access control are thus required to be more innovative, and utilizing minimal information. In the study of the protocol schemes that follows in this work, the problems of modulation, synchronization, coding and the like are all assumed to be solved. The main concern is the mechanism responsible for sharing the channel resource.

The ability to utilize the vast transmission bandwidths of optical fibers rests in designing communication protocols that facilitate highly coordinated signal multiplexing. The design of protocols for high speed optical networks thus necessitates a requirement beyond those typically demanded of a multiaccess scheme. In this type of network, there is need for a frequency (subchannel) assignment to pairs of communicating stations as well as the multiaccess control of each of the subchannels operating at the particular frequencies. The need to coordinate this two-dimensional allotment of frequency and maintenance of multiaccess control, subject to the constraints of efficient traffic throughput and delay demands, adds a significant degree of complexity to the design of the protocols. A useful visualization of the problem is the simultaneous coordination and control of multiple, parallel multiaccess channels, for which the end users must still be totally interconnected.

## 1.4 Thesis Outline and Organization

The aim of this study is to propose a protocol scheme for the control of communication on future optical networks based on a better understanding of multiple access methods in the single satellite channel environment. This study thus seeks to provide a clear classification and categorizing of satellite protocol methods and an understanding of their performance capabilities under various operating conditions. Together with an examination of the work presently done in the area optical network protocols, the understanding gained from the study of satellite protocol designs will be used as a basis for proposing a protocol approach for optical network communication. The relevance of this approach in studying satellite schemes is further supported by the recognition of the characteristics of satellite multiaccess protocols that have already appeared in the literature of optical network control methods.

A major component of the thesis work, therefore, is the compilation and examination of an exhaustive list of protocol schemes developed in the satellite communication arena. An analysis of the foundations of satellite protocol approaches and the respective performances achievable will provide the basis for a scheme that can address the demands of parallel, cross-linked, multiaccess channels. This optical network protocol will thus build from the usefulness of satellite multiaccess schemes while incorporating mechanisms to handle the differing realities of the future optical network.

Chapter 2 provides a basic introduction to satellite communication and the general types of schemes that have been proposed in that context. While many proposals exist for multiaccess communication, this study focuses on methods that are applicable in environments of long propagation times relative to transmission times. In discussing the basics of multiaccess satellite communication a general framework is presented establishing a range of categorizes into which existing protocol schemes can be ordered. This provides a context for the examination of the schemes which then follows.

Before the compilation of satellite protocols is examined, however, some discussion is provided in Chapter 3 on the various requirements of a protocol and the factors used in assessing operating performance. This is followed by a presented summary of an extensive list of protocol schemes. The schemes are ordered according to their general characteristics and the categories that they most closely represent. This side-by-side examination allows for the establishment of a clearer picture of the similarities and subtleties between differently proposed methods. It better allows for a conclusion to be derived on the extent to which protocol bases exist.

Chapter 4 seeks to extend the conclusions of the work done in Chapter 3 by examining a number of recently proposed optical communication schemes in an attempt to foundations of the various protocol methods carry through. This would provide addition validation of the premise of this thesis work - namely that by better understanding and establishing a foundation from which protocols have been derived, we can have a clearer process in the proposing of future schemes. This therefore is the preparation for the presentation and analysis of a new optical communication scheme which follows.

The work of the initial four chapters has defined more clearly the nature of the problem being addressed in the development of an appropriate mechanism for achieving access control in the multi-channel optical environment. The details of the protocol scheme introduced, though new in its combination of elements, represents a bring together a number of the ideas encountered in our earlier study. Their application, however, to the realities of optical network communication makes the contribution unique. Chapter 5 thus explains the operational details of the method proposed and provides the necessary justifications of the included features. The qualitative analysis and implications of the scheme's operation is also presented.

In chapter 6 we place into a larger context the undertaken technical study which has led to the proposal of a new operational multiaccess scheme. This larger framework involves an understanding of the process of technological innovation and development. The approach adopted in this study was

based on an essentially technical rationale. How this has affected the way in which we proceeded and the results we obtained is of importance in future technical studies. We thus seek to perform a retrospective examination of the process involved in the work undertaken. From this analysis we can identify potential key determinants of technological development in the case of technical protocol design. This leads to the defining of technology policy considerations for other similar studies.

Chapter 7 is the final brief overview of the results our work; technical and non-technical. We conclude by providing a brief summary and conclusion of the technology study performed, the technical results obtained and the issues raised by the process employed. Thus, the outcome of this study has not only being the derivation of a useful technical result, but also a broader understanding of the factors that may shape the development of similar technologies.

# Chapter 2

# Multiaccess Satellite Communication

## 2.1 Introduction to Satellite Communication

A satellite and its component transponders in a geostationary orbit above the earth provides long-haul communication capabilities. It can receive signals from any earth station in its coverage pattern and can transmit signals to all such earth stations (neglecting issues of spot beam usage). Full connectivity and multi-destination addressing can both be readily accommodated. In the most general sense, the communication satellite network functions simultaneously both as a local network supporting directly attached hosts and terminals of differing capabilities and requirements and as a transit network with attached internetwork gateways. That is, the satellite network may also support the connection of other smaller networks.

The focus of our concern throughout this study, however, will be on the local networking capability of the system's operation, with gateways and other interconnecting modules viewed simply as additional nodes of the local network. The issues of network interfacing, though also intrinsically important to overall communication, represents an independent element of the system that must be addressed in its own right. It is the operation of the central local network that forms the core of the satellite communication system.

21

The satellite can be viewed as a centralized power and bandwidth resource which must be shared among a number of communication links. These links have several important characteristics form the user point of view: data rate, error rate and most importantly propagation delay. Link data rates can extend to hundreds of megabits per second depending on the application and the choice of system parameters. For fixed-power and bandwidth parameters it is possible to tradeoff between the link data rate and error performance; decreases in the data rate typically permit better error performance. The most important characteristic relevant to this study is the long propagation delay for a single hop. This round trip propagation delay (to and from a geosynchronous satellite) of approximately 0.27 s which is usually long compared to the transmission time of a packet, has a major impact on bandwidth allocation techniques and on error and flow control protocols.

Another important characteristic of a satellite communication system is its coverage. This area, determined by the transmit and receive antennas of the satellite, is the region over which stations can access the satellite and from which they can receive its transmissions[3]. It is a main determinant of the system as the satellite can provide broadcast capability at any time to all stations in its coverage area. Due to the broadcast nature of the communication channel, each station (user) is able to listen to and receive any messages transmitted by any other terminal in the network, including itself. The combination of multiple access by dispersed users and the broadcast capability allows stations to be formed into a fully connected one-hop network.

Telephone users make their calls at random. So do computer terminal users. A central problem of telecommunications is that of assigning communication channels in a sufficiently flexible manner that users can have access whenever they have a whim to communicate. Broadcast satellite communication thus faces the same type of challenges as other communication networks. Namely, the need to allocate

---

[3]    The use of Space Division Multiple Access (SDMA) methods which are based on the separation of users due to different satellite coverage areas, is not of relevant concern in this study. For our purposes, stations in the network are connected in a single hop to the satellite.

a limited transmission capacity among a number of users whose demands are varying over time. As the number of users increase, and the variability of their demand becomes more random, so too must the sophistication of the allocation approaches improve in order to maintain overall system performance and resource utilization efficiency.

## 2.2 Multiaccess Communication

Multiaccess protocols thus arise out of this need to share the common communication resource among a number of independent users. Some means of channelization is required such that transmission capacity can be provided to active users. To satisfy the demand requirements of users, the communication resource could be allocated by way of a i) fixed assignment, or on a ii) variable or demand assignment basis. With demand assignment an additional channel must be set aside for coordination signalling among users. (Access to this signaling channel itself presents a multiaccess problem). Fixed assignment divides the transmission bandwidth so that independent of its activity, each station in the connected network is provided a fixed allocation of the resource for its use.

The two common forms of fixed allocation are based on the orthogonal schemes of: synchronous time division multiple access (TDMA), allows users fixed time intervals during which they can transmit using the entire communicating bandwidth, and frequency division (FDMA) where the entire bandwidth in divided into a number of separate communication subchannels for use by individual network members. In the context of fixed assignment, FDMA suffers from the disadvantages of wasted bandwidth due to channel separation, lack of flexibility in dynamic allotment of bandwidth, and lack of broadcast operation. The major disadvantage of TDMA is the need to provide rapid burst synchronization and sufficient burst separation to avoid time overlap.

Code Division (CDMA) is also another means of fixed assignment allocation in which simultaneous channel transmissions are differentiated on the basis code signatures. These signatures could be defined in the time domain or in the frequency domain. Spread Spectrum is the form of CDMA whereby each user is assigned a particular code sequence which is modulated on the carrier with the digital data modulated on top of that. Spread spectrum distributes the signal over a wide band of frequencies and then compresses it at a receiver that identifies the signal by a code. In such CDMA schemes a number of users occupy the entire transmission bandwidth all of the time. Their signals are encoded so that information from an individual transmitter can be detected and recovered only by a properly synchronized receiving station that knows the code being used. Because the power of a SS signal is so diffuse it provides little interference with other signals sharing the same bandwidth. The appeal of SSMA lies in its potential for uncoordinated access and its security and interference rejection characteristics. It however requires a large processing gain which limits its attractiveness for networks with large numbers of users. Spread spectrum communication has nonetheless become a complete area of study in its own right but is not further address in this paper.

Because these fixed assignment schemes make allocations independent of the changing needs of the various users they tend to be inefficient in their use of the central communication resource except in situations where the traffic being generated by network stations are evenly distributed and regular over time. Because of the more general need to adapt to variations in traffic level and statistics, the more effective approach lies in demand assignment schemes where the communication bandwidth is dynamically allocated on this basis of the changing environment conditions.

Another main property of a multiaccess protocol is its ability to as much as possible provide immediate access. That is, as the traffic load in a network goes to zero, the delay that a user experiences before transmitting a packet should accordingly fall. In schemes where there is minimal coordination overhead prior to transmission, the immediate access property can be more closely realized. However, where control information interchange is required prior to the transmission of packets, the minimum

delay involves at least a round trip propagation time (which in satellite networks is not negligible). In the performance discussions that follow later in this work, it is seen that the desire of immediate access is traded-off for improved utilization of the central communication resource.

A number of different approaches to demand assignment have been developed for satellite communication. These can be grouped into those designed to handle circuit-switched voice traffic, and those designed primarily for packet-switched data traffic. As shall be seen, circuit-switched approaches are not very well compatible with the demands of data applications and interactive communication transmissions usually characterized as bursty. This is because of the particular characteristic of voice oriented traffic - essentially, the long duration time of calls compared to the time required to make new circuit allocations - which is better accommodated by the setting up of new transmission paths on each occasion that a station wishes to transmit. Packet communication on the other hand is based on the idea that part or all of the available resources are allocated to one user at a time but for just a short period of time. This better allows for the interleaving of short duration transmission requests.

The emphasis of this study will be to examine packet-switched approaches since they not only address the demands of future optical networks, but also in general better represent present communication needs. The packet-switched approaches to demand assignment of the communication resource can be separated into two main categories: random access or contention, and regulated access or reservation schemes. The regulated access approach can be further divided into implicit and explicit reservation schemes, as shall be defined below. There is also a class of schemes that seek to provide dynamically adaptive capability between the alternatives of random and regulated access, or which employ a mix of these approaches.

Access control methodologies can also be distinguished on the basis of the way in which the mechanism for maintaining control over channel allocations operates. That is, whether allocations are performed

in a distributed manner with all active stations responsible for channel assignment, or in a centrally controlled manner where a single station (with some backup) coordinates assignment. In the satellite communication environment the consequence of large propagation delays restricts the use of centrally controlled schemes. However, as schedule processing is taken on-board the satellite, the delay hops associated with centralized control is reduced. We thus mainly focus on distributedly controlled methods except in cases where some superior performance advantage can be gained from the use of a centrally controlled scheme.

## 2.3 Circuit-Switched Approaches

Two of the early systems designed for circuit-switched voice traffic are SPADE [EW 72] and MAT-1 [Sch 69]. In the SPADE system, transponder capacity is divided into subchannels using FDMA with all except one of these subchannels dynamically assigned to requesting stations. The remaining subchannel, a wider bandwidth channel is used as a TDMA order wire, with each station permanently assigned a one time slot per transmission frame. This common signaling channel allows stations to make requests for voice channels from the available pool. A station needing to transmit examines the table of available frequencies and then seizes a free channel. Control is decentralized. No specific stations are in charge of allocating channels. It is all handled across the common signaling channel.

In the MAT-1 approach all circuits are derived from the transponder using TDMA channelization. Each subchannel consists of a time slot defined to contain 8 bits of PCM digital voice data, with a frame size of 125 microseconds. The subchannels within the frame are partitioned into groups, with each station assigned a group. The number of subchannels in each group is reallocated periodically so that stations with heavy demand can have a larger share of the total.

Both the SPADE and the MAT-1 systems are efficient for the voice calls they were designed to handle. If these circuits were used for bursty traffic, however, these systems would not appreciably improve channel efficiency over the case of no demand assignment. A great deal of overhead would be spent in the constant setting up of new circuit paths. New techniques were thus derived to provide demand assignment satellite capacity for packet transmissions.

## 2.4 Packet-Switched Approaches

### 2.41 Random access:

· This approach is characterized by the ALOHA method originally developed for ground based radio [Abr 73]. Subject to constraints designed to maintain system stability, each station which has a packet to send simply sends it. If this packet overlaps in time with transmissions from one or more other stations, a conflict occurs and the reception is garbled. The broadcast property of the channel usually allows each station sending packets to monitor the success or failure of its transmission and to queue unsuccessful packets for retransmission. This retransmission can take place after a period of time equivalent to the maximum station-to-station propagation-plus-transmission time. In this scheme, no central control or coordination of stations is required. When a packet is to be transmitted the chance is taken to transmit with a resulting possibility that a collision or conflict occurs.

Because of the variability of station traffic and the possibility that transmissions and conflict-retransmission may lead to runaway conditions, it is necessary to augment this free transmission protocol with some further control measure. This could involve varying the probability of immediate transmission of packets on arrival at a station as well as packet retransmissioms following a collision. For optimum system performance, this additional stability control must be dynamically able tot respond to the traffic state and loading of the network. Under low load conditions uninterrupted transmissions should be encouraged, while for high load and large numbers of queued station packets,

the number of retransmissions should be appropriately regulated. In general however, random access schemes are characterized by their simplicity of implementation.

Carrier sensing multiple access (CSMA) methods can also provide improved random access performance (through reduced transmission collisions) by listening to the channel before beginning a transmission. This is only viable, however, in situations where the propagation delay between any source-destination pair is small compared to the packet transmission time. As such, these schemes are not applicable in a satellite communication environment. Similarly, approaches based on global scheduling are inefficient in the distributed, high propagation delay environment.

### . 2.42 Regulated Access:

To improve the throughput utilization of the communication channel a reservation approach is used to the transmission of packets. The objective of reservation protocols is to avoid collisions entirely. In these schemes, each station sends a reservation packet containing information regarding channel time required by the station. These requests are received by the the central controller (in a centrally controlled system), or by each active station in a distributed control scheme, and are subsequently used to assign channel service times to the requesting stations, according to the system's access-control discipline.

Channel bandwidth resources are thus appropriately divided into reservation and service components. The first component is used to process reservation request packets while the second is dedicated to the actual transmission of data or message packets. This division of channel bandwidth can be done by either of time or frequency multiplexing. An advantage of time multiplexed channels, however, is that the partition may be variable. With a variable partition the service component fraction of the channel bandwidth can be made very close to one under conditions of heavy load. With a fixed bandwidth partition there is a fixed capacity loss. The key problems to be addressed in these reservation schemes are: i) the implementation of the reservation subchannel, and ii) the implementation of a global queue.

The implementation of a global queue becomes more of a significant concern in the distributed satellite environments which we will examine.

## 2.421 Implicit Reservation:

In protocol schemes of this type, reserved access to the communication channel is based on a reservation-by-use approach. Initial access is gained through contention (random access) among users. These protocols utilize a frame format in which the communication medium is divided into discrete transmission intervals with corresponding sub-intervals for individual station transmissions. As in all that employ a frame format, the time duration of the frame must be at least one round trip propagation time to allow reservation information to be conveyed to all stations before each schedules it own transmissions. When a station uses a contention slot successfully, the transmission slot is assigned to that station until it stops using it. This is thus an approach which mixes contention with reserved assignment.

This implicit reservation allows stations with high traffic rates to have greater use of the communication channel by capturing more transmission slots for their exclusive use. Control is distributed in that each station executes an identical assignment algorithm based on global information available from the channel. Stations are therefore required to maintain a history of the usage of each channel slot for one frame duration.

## 2.422 Explicit Reservation:

In this demand assignment approach a portion of the channel resource is used (as an order wire) by stations to explicitly make a reservations for transmission of a packet or packets. The reservations may be sent in separate subframe(s) distinct from the frames in which messages are transmitted, or they may be combined with message transmissions ("piggy-backed") or both. This division of the channel into reservation and service components could be made on both a temporal basis or through a frequency

separation. In addition, the assignment of message transmission times resulting from reservation may be made centrally by one station, distributedly by all stations , or by a mixture of both techniques.

Since the making of reservations involves a situation similar to that encountered by stations in obtaining access to the main transmission channel, various approaches can be used. The multiaccess problem on the reservation channel could again be solved using fixed or demand assignment schemes. Thus in schemes where the assignment of the reservation channel is not fixed the ratio of reservation subframe sizes can be varied according to traffic loads, with the entire transmission frame used for making reservations when there are no messages for transmission. As shall be explored in the course of this study a number of different protocol approaches hinge around the choice of methods used to make · packet transmission reservations. In all types of schemes, however, it is important to keep the segment of the transmission bandwidth dedicated to reservation as limited as possible. This avoids unduly occupying channel capacity that could be used for data message transmissions.

## 2.43 Mixed Protocol Schemes

Another category of protocol approach to the multiaccess communication problem is that involving mixed schemes. Unlike pure random access or pure reservation strategies, or even pure TDMA, which possess particular operating strengths but within certain inherent limitations, these protocols are distinguished by their integration of several different access techniques in the same system. The objective is to take advantage of the particular attributes of the different accessing scheme by having the modes coexistent and operating simultaneously in a manner which derives a potential benefit. Mixed schemes represent a progression in the move away from methods whose effectiveness is confined to narrow operating ranges.

### 2.44 Adaptive Protocol Schemes

As shall be seen in the next chapter, there have been quite a number of methods proposed for dealing with the multiaccess communication problem. Different classes of schemes tend to have their own advantages and limitations. With the exception of perfect scheduling, which is not achievable in a distributed environment, none of the protocol methods can boast superior performance over the entire range of system requirements. If a scheme is able to perform nearly as well as perfect scheduling under low input traffic conditions, then it tends to be limited by achievable channel capacity. On the other hand, schemes that are efficient under conditions of high system utilization tend to have prohibitive access control overhead requirements when utilization is low.

The resolution of this tradeoff effect in protocol operation can be achieved to some extent by applying dynamic measures to specific access schemes. Such mechanisms will allowing for improved performance may not be able to entirely overcome the limitations of an access scheme. For example, while dynamic adjustment can provide flexible reservation schemes, it may not be able to change the access control to completely random access when new traffic condition may deem such an approach more efficient. Rather than dynamic adjustment, therefore, what may be more adequate is a strategy for choosing an access mode which is itself adaptive to the varying traffic needs. In such a protocol approach, among the crucial factors would be the type of information required by the adaptive strategy and the implementation of the information acquisition mechanisms.

## 2.5 Establishing Performance Criterion

Data communications have very diverse requirements, ranging from inquiry-response systems with intermittent traffic to file transfers with large volumes of data. In addition, some of these data transmissions have user-specified delay constraints that need to be met. Access control protocols are

required to to be responsive to these varied demands as efficiency in channel resource utilization comes from the ability to multiplex different traffic types. An appropriate measure of this traffic carrying capacity of a channel is the aggregate throughput rate in number of messages (or packets or bits) that can be transported per unit time, while satisfying specified delay constraints. Some measure of a protocol's performance can thus be determined from its channel throughput versus delay tradeoff characteristic. This performance characteristic seeks to ensure high bandwidth utilization and low message delays.

A number of other attributes are, however, also important. The ability for an access protocol to simultaneously support traffic of different types, different priorities, with variable message lengths, and differing delay constraints, is also essential if diverse traffic types are to be supported. These requirements nonetheless need to be achieved at minimum overhead cost and with minimal imposed time penalties if high throughput and low delays are to be maintained.

If we let $C$ be the channel transmission rate (capacity) in packets (or bits) per second and B, the number of packets (or bits) successfully transmitted per second, then the throughput of the channel, $S$, can be defined as the ratio of the rate of successfully transmitted packets to the transmission rate of the channel. That is, $S=B/C$. Thus channel throughput is a normalized quantity between 0 and 1. Delay on the other hand is defined as the time between the arrival of a packet at a station and the time of the packet's successfully transmission. Packet collisions (unsuccessful packet transmissions due to interference caused by a whole or partial overlap of transmissions from competing stations) thus increase the average delay time associated with a particular protocol scheme. So too does the time a packet must remain queued at a station before transmission is attempted.

The performance of a multiaccess protocol is strongly dependent upon the input traffic model (packet arrival at each station) and the network loading. This is also affected by the use of flow control mechanisms as well as different buffering conditions. However, since a number of performance

evaluations of multiaccess protocols have been separately conducted, the assumptions and conditions used have not been identical. This makes performance comparisons less definite. Another limitation stems from the difficulty in exact performance analysis and the need to make certain assumptions so that the problems become tractable. For example, the assumptions of buffered or unbuffered users, the choices of input source models and the assumption of statistical equilibrium. These factors therefore prevent very precise comparisons of different protocol schemes when performance has not been derived on an identical basis. The subtleties of the different evaluation conditions must also be considered.

Tasaka [Tas 84] presents the use of an *equilibrium point analysis (EPA)*,[4] an analytical technique which provides the means for deriving the performance of different protocol schemes under identical dynamic system conditions. Within the scope of our work, however, this technique will not be applied to the performance evaluations of the various protocol schemes examined. With analyses already performed by a number of researchers, though not using EPA, there is still sufficient information to compare a wide variety of schemes since each analysis usually makes comparisons with existing methods as well as the ALOHA base reference. As a result, we will rely on the evaluations presented in the individual papers and use this to cross compare among the entire list of examined protocol methods.

It is instructive to understand how different condition factors affect performance. A good example concerns the stability problem that affects contention-based protocol schemes - the fact that the system could become overwhelmed with excessive collisions and the need for retransmissions such that throughput falls to zero. The performance analyses usually consider situations of statistical equilibrium neglecting the potential instability that results under certain loading conditions. In reality, the performances derived for the system under contention-based protocols is only achievable for a finite period of time [Tas 84]. Performance comparisons thus tend to be made between static models.

---

4    The equilibrium point analysis (EPA) is a technique that allows multidimensional Markov chains to be easy analyzed. Dynamic behavior of a satellite communication system can thus be more readily evaluated using Markovian models.

A truer comparison of these protocols would have to be conducted under conditions which guarantee system stability (as done by Tasaka).

Another issue of note in our consideration of protocol methods is the difference in performance that results from evaluations based on different message lengths. That is, between messages of very few packets and those containing a greater number. In some schemes the particular mix of message lengths is also an important performance determinant. These differences must thus also be considered when comparing the independently derived performances of alternate schemes.

Though not directly addressed in most protocol schemes, automatic repeat request (ARQ) mechanisms are necessary in practical communication protocols for dealing with multi-packet messages. This stems from the fact that data sources generally produce messages which vary widely in length. In order to accommodate these diverse lengths in fixed time slot synchronous protocols, messages must be subpacketized. The alternative is the use of large length slots which result in wasted channel capacity when short messages are transmitted. The inclusion of these ARQ techniques results in some measure of degradation in the nominal performance of a protocol scheme which has been evaluated without the consideration of message packets assembly, ordering etc.

Jacobs et al [JBH 78], in defining the requirements of a general purpose satellite network, provided some useful insight into various qualities of satellite protocol schemes that may make them more versatile and so contribute to increasing their performance value beyond that of the throughput versus delay tradeoff. While throughput versus delay still remains our central performance measure, the distinguishing attributes of some protocol schemes allow them to be judged highly in categories where other performance criteria are also of dominant consideration. For example, a number of schemes though efficient in bandwidth utilization and message delay are unable to work in communication environments requiring the ability to handle prioritized traffic.

By considering some of these additional qualities prescribed by Jacobs, we can have a more complete framework for the protocol examination to be undertaken through our ability to recognize the value of certain attributes provided by different schemes. Some of the more important requirements include: i) satisfaction of multiple delay constraints, ii) multiple priority levels, iii) variable message lengths, iv) stream traffic, v) fairness and vi) robustness. In the examination and analysis of protocol schemes that follows, we will thus maintain our focus on the criteria of throughput versus delay but will nonetheless be aware of the limitations of direct numerical comparisons, as well as recognize the other attributes that may make certain schemes very attractive in specially defined operating environments.

In the protocol examination to follow, we shall look closely at these demand assignment foundations presented above and the way in which they are implemented in achieving access control. It will be important to understand the strengths of the different approaches and the success each achieves under different operating conditions. The descriptions of the various protocols are meant to focus on the essential details of the schemes as well as provide an understanding of their operation[5]. Unless particularly relevant to the unique functioning of a particular scheme, the issues of practical implementation are neglected in the presented descriptions. Performance measures are provided and where possible comparative performance data given on different schemes.

---

[5]    The broadcast nature of the satellite communication system allows feedback on a given transmission to be received by all stations after one roundtrip propagation delay. In the protocol schemes below, the feedback condition generally assumed is that of detection of 0, 1 packet or a collision. Due to topographical and environmental conditions the station receivers are prone to certain errors, including noise errors. The occurrence of such non-ideal events are however neglected and so too are the contribution of capture or erasure effects.

# Chapter 3

# Survey of Proposed Satellite Communication Protocols

## 3.1 Random Access Protocol Schemes

These protocol schemes are characterized by their minimal coordination requirement for the transmission of data packets on the communication channel. Access is not controlled but determined by the random generation of packets at network stations. As a result, a key element of these schemes is their simplicity of implementation.

3.101) *ALOHA [Abr 73]*

This protocol strategy represents the foundation from which multiaccess communication in the satellite environment was developed. In this protocol, each network station upon receiving a new packet, transmits it immediately. Collisions therefore inherently occur as no attempt is made to regulate access to the communication channel. Collided packets are then retransmitted after some randomized delay. In contrast to fixed assignment approaches, small delays are achieved if collisions are rare. This has to be traded-off, however, against the potentially large delays that may occur when there is heavy network traffic and also the instability that may result from excessive collisions and perpetual retransmissions.

This protocol method while having the advantages of ease and simplicity in implementation, is however, only able to achieve a maximum channel utilization of about 0.18 (1/2e). This effectiveness is also only limited to bursty transmission environments. The familiar performance characteristic (min delay at low traffic throughput rising to infinite delay as throughput approaches the 1/2e limit) associated with this protocol is based on assumptions of: a) statistical independence of channel traffic, b) infinite user population, and c) statistical equilibrium. As detailed in [Lam 79], these assumptions have been investigated and while the protocol is potentially unstable without adaptive control, the results based on these assumptions are still robust under less extreme conditions.

3.102) *Slotted ALOHA*

The Slotted ALOHA protocol is very much the same as ALOHA with the additional requirement that the channel is slotted in time. Users are required to synchronize their transmissions into fixed length channel time slots. The limiting of transmissions to time slot boundaries results in the slotted ALOHA having a maximum channel throughput that is twice that of the unslotted case (about 0.36). It was also shown [Abr 73] that when the traffic distribution is unbalanced with a mixture of high rate as well as low rate users, the maximum throughput of this protocol could be considerable improved. The performance characteristic is similar to that of the pure ALOHA. This scheme however has the additional requirement of synchronization of users.

3.103) *Diversity ALOHA [CR 83]*

In this random access scheme slotted ALOHA is employed with the additional element that multiple copies of a given packet are transmitted. In the frequency diversity version, packets are simultaneously transmitted on different channels. In the time diversity scheme copies of the packets are transmitted on the same channel but spaced apart by random time intervals. In time diversity, two approaches can be employed involving the transmission of either a fixed number or a random number of copies of a packet. This scheme reduces average packet delay by minimizing the need to wait an entire round trip

propagation delay before retransmitting a collided packet. It is assumed that some arrangement exists at the receivers to reject all copies of a packet once a successful reception has been achieved.

It is found that multiple transmissions give better performance if throughput is somewhat below its maximum. Also, if the probability that a packet fails a certain number of times or more is not to exceed a certain limit, the multiple transmissions usually give a greater throughput.

### 3.104) *Stabilized ALOHA*

Various means exist for the stabilization of the ALOHA random access operation. The approaches can be based either on some measure of regulation in the initial transmission of packets, regulation of retransmissions, or some combination that may also include feedback control. Essentially, when idle slots occur, the packet transmission/retransmission probabilities are increased. They are however decreased when a collision occurs. One stabilizing approach is Rivest's [Riv 85] pseudo-Bayesian algorithm. The scheme differs from slotted ALOHA in that newly arrived packets, rather than being immediately transmitted, are considered backlogged and thus transmitted in the next slot with some probability just as packets involved in previous collisions. The determination of the retransmission probability is based on an estimate of the number of backlogged at the beginning of each transmission slot.

### 3.105) *Selective Reject (SREJ) ALOHA [Ray 84, Ray 87]*

The Selective Reject ALOHA is a technique that incorporates a particular automatic repeat request (ARQ) retransmission method for dealing with the transmission of multi-packet messages. It has been included here because of the interesting result obtained.

Unlike some other ARQ strategies that retransmit either the entire message or all packets starting at the first erroneous packet, the SREJ protocol only retransmits those packets which were received in error. Retransmission of entire messages is thus avoided when only partial overlaps occur. This

however requires that packets be numbered and addressed, and receivers be capable of buffering and reordering of message packets. This implies slightly greater overhead and additional equipment complexity. It was shown [Ray 84] that the throughput of SREJ ALOHA (which retains the advantages of asynchronous operation) is identical to that of conventional single-packet slotted ALOHA[6], irrespective of the message length distribution. This scheme thus leads to significant capacity as well as stability improvements over conventional multi-packet message ALOHA.

.

---

[6]     Given subpacket overhead however, the maximum achievable throughput of practical applications is in the region of 0.2 - 0.3 [Ray 87].

## 3.2 Reservation Protocol Schemes

These protocols are based primarily on the reservation of the communication channel by individual stations wishing to transmit. Some segment of the channel is thus maintained for reservations. The differentiating characteristic in a number of these reservation schemes is the manner in which the access to the reservation channel is controlled and the nature of the allocation and implementation of the reservation channel. A desired objective, nonetheless, is that a minimum overhead is incurred in the making of reservations.

### 3.21 Implicit Reservation Schemes

3.211) *Reservation-ALOHA [CRW 73]*

This protocol represents the implicit reservation scheme. It is based on a slotted time axis, where the slots are organized into frames of fixed size. The frames must be of length greater than the satellite propagation delay. Packets arriving at a station are placed in a queue. Unused slots (those free or containing a collision in a previous frame) are then accessed by all users in a slotted ALOHA contention mode. This contention is however regulated to some extent by having the probability of a packet transmission at a station varied according to the length of the packet queue at that node. A user who has successfully accessed a slot in a given frame is then guaranteed access to the same slot in the succeeding frame and this continues until the user stops using the slot. The users maintain the history of usage of each slot for one frame duration.

This type of reservation by use approach is only effective if users generate stream traffic or long multipacket messages. Performance will degrade significantly with single packet messages, as every time a packet is successful the corresponding slot in the following frame is likely to remain empty.

## 3.22 Explicit Reservation Schemes

### 3.221) *First-in First-out (FIFO) Reservation Scheme [Rob 73]*

In this scheme, channel time is slotted with one slot (after every $M$ slots) being periodically divided into a number of minislots. These minislots are used for the transmission of reservation packets on a slotted ALOHA contention basis. The remaining (data) slots are then used for reserved transmissions. When a data packet or a multi-packet block arrives at a station it transmits a reservation in a randomly selected minislot of the next minislot group. Once successful, the number of transmission slots requested is seen by all stations (after a propagation delay) and added to the count of reserved packet slots. There is thus a common reservation queue for all stations and by broadcasting reservations they can claim a space on the queue. If the reservation attempt is unsuccessful, the station retransmits the request in the next minislot group following the determination of the request collision.

There are two operating states, ALOHA (for making reservations) and Reserved (for actual data transmissions). If the reservation queue is empty the channel operates entirely in the ALOHA state to accept reservation requests. Each station is required to maintain information on the number of outstanding reservations (the queue) and the slots at which its own reservation begins. These are determined by the FIFO discipline based on the successful request packets received. Queue status information is also transmitted in data packets to assist synchronization.

The FIFO reservation scheme, operated with either a TDMA or a slotted ALOHA reservation subchannel, offers delay improvements over fixed assignment TDMA [Lam 77]. When compared to ALOHA, a higher system capacity is achieved but at the expense of a higher delay at low channel throughputs (due to higher overhead). Tasaka [Tas 84] using his equilibrium point analysis showed that the use of ALOHA on the reservation subchannel, for a fixed frame length system (differing somewhat from the variable frame length FIFO scheme), gives better delay performance than the use of TDMA, though at the expense of a loss in maximum system throughput for single packet messages. As the number of packets per message is increased, however, ALOHA reservation is able to achieve a

greater throughput while still preserving its lower delay advantage. The figure below [Tas 84] shows the performance comparison of fixed assignment TDMA (TDMA), Slotted ALOHA with go-back-N automatic repeat request (S-ALOHA /GBN), Reservation ALOHA (R-ALOHA), modified FIFO (fixed frame length) with ALOHA reservation (ALOHA-Reservation), and modified FIFO with TDMA reservation (TDMA-Reservation), for a system of 100 users (M) and messages whose lengths are geometrically distributed with a mean of 10 packets (h).



Figure 3.221    Delay-throughput performance comparison of TDMA, S-ALOHA/GBN, R-ALOHA, modified FIFO (ALOHA reservation) and modified FIFO (TDMA reservation). Source [Tas 84] p. 1579.

3.222) *Round-Robin (RR) Reservation Scheme [Bin 75]*

This protocol is a modification of conventional fixed assignment TDMA in which provision is made for the dynamic assignment of unused channel time slots. The scheme employs a frame format with the frame length of duration greater than (or equal to) the propagation delay time and such that the number of slots in a frame is at least as large as the number of network stations. Each station is permanently assigned a slot within a frame which it must use for its transmissions. In each transmission frame, slots not needed by their owners are allocated, on a round-robin (RR) basis, to stations with additional packet transmission requirements.

This dynamic assignment of unused slots is achieved through the use of a reservation system and the operation of a distributed network queue. The reservation information is sent as part of the overhead of the data packet in the owned slot of a active station. These reservation requests are then collected by all stations. Each station thus maintains a queue table which contains an entry for each node (which also indicates inactive nodes) and allows for the RR assignment of available slots.

A station recovers its own slot which was being used by deliberately causing a conflict in that slot. The result of this event received one round trip later by all nodes indicates that a node wishes to regain its slot and so the slot is removed from dynamic assignment. The periodic transmission of a station's queue table is employed to assist synchronization in this distributed scheme. Other than round-robin assignment of unused slots can be employed in this protocol approach. The figures below provide performance comparisons of the RR scheme for long (multi-packet) messages with that of TDMA.

Binder's [Bin 75] results also shows that the performance of this scheme is unsatisfactory for short messages, with average delay being even greater than that of TDMA for almost the entire throughput range. This is possibly accounted for by unused reservation overhead as well as recovery overhead penalty incurred when stations wish to regain use of their assigned slots. The scheme is also noted to give significant performance gains as the traffic imbalance among the nodes increases.

Comparing Binder's results with those shown above [Tas 84] (section 3.221), it appears that the RR scheme does not give the same extent of performance gain over TDMA as that of Reservation ALOHA or the modified FIFO reservation scheme. It should be pointed out, however, that we have used Tasaka's results for the performance of Reservation ALOHA and the modified Robert's FIFO scheme rather than the results originally presented by Binder. This is justified on the basis of Tasaka's methods, the recency of his data, and the fact that Binder's comparative results for Reservation ALOHA were preliminary simulation findings.

Note: In the figure below the vertical scale is specified in terms of round trips (RT's), where RT = 0.27s.

Figure 2. Long Message
Packet Delay vs. Thruput

**Figure 3.222** Delay-throughput performance comparison of TDMA, Roberts FIFO, Reservation ALOHA,, and Binder's Round-Robin reservation for long messages. Source [Bin 75] p. 41-3.

3.223) *Split-Channel Reservation Multiple Access (SRMA) [TK 76]*

In this protocol method, available channel bandwidth is either time divided or frequency divided between reservation and actual message data transmission. One channel is used to transmit control information, the second used for the data messages themselves. With this configuration there are many operational modes based on the access method used on the reservation subchannel.

In the request/answer-to-request/message scheme (RAM), the bandwidth allocated for control is further divided into two channels: the request channel and the answer-to-request channel. In this scheme, the request channel is operated in a random access mode. A central station performs the scheduling and issues on the answer-to-request channel the time at which a requesting station can begin transmissions.

44

Another version of this protocol request/message (RM) involves having just two channels: the request (control) channel and the message channel. Requests received by the scheduler are queued and serviced according to some appropriate algorithm (e.g. "first-come first served" or some other algorithm which allows for priority). When the message channel is idle an answer-to-request is sent out on the channel with the ID of the station that should begin transmission. The station recognizing its ID starts transmitting on the message channel. Requests that remain unanswered beyond a given interval of time are reissued. The results of the protocol versions are shown below. The basic approach here is similar to that of Robert's FIFO scheme (though the implementation and control on the reservation channel may differ). This scheme is representative of the basic multiaccess reservation approach.

In this protocol an optimal bandwidth allocation exists between control and message channels. The minimum packet delay given in the figures refer to the delay obtained when the system operates at this optimum.

Figure 3.223a    Delay-throughput performance of SRMA with ALOHA reservations for different control information overhead values, η. Source [TK 76] p. 843.

Figure 3.223b    Delay-throughput of ALOHA, Slotted ALOHA and Perfect Scheduling under the same conditions as SRMA in Fig. 3.223a. [TK 76] p. 844.

The delay is given in units of number of bits per message packet, $b$, divided by the total channel bandwidth, $W$. $\eta$ is a measure of the overhead due to control information (ratio of bits per control packet to bits per message packet).

3.224) *Priority-Oriented Demand Assignment (PODA) [JBH 78]*

This protocol represents a outgrowth of the FIFO scheme (see section 3.221 above) with the addition of more sophisticated scheduling and operating features. Channel time is divided into two basic subframes, an information subframe and a control subframe. The information subframe is used for the transmission of packets and stream traffic that have been scheduled. These transmissions also contain contain control information such as acknowledgements and reservations in their headers for further reservations by the transmitting station. The control subframe is used to send reservations that cannot be sent in the information subframe in a timely manner. The scheme thus incorporates both implicit (through headers of transmitted packets) and explicit reservations. The control subframe must be used for initial reservation access by a station that has no impending transmissions (it can also be used for priority scheduling). The information subframe is further divided into segments that allow for the integration of centralized and distributed assignments.

Access to the control subframe (which together with the use of packet headers constitutes the reservation channel) is dependent on the system environment. For a small number of users, a fixed assignment access approach may be used (referred to as FPODA system operation). For a large number of low-duty cycle users, mixed users, or unknown traffic requirements, a random access or contention alternative may be used, such as slotted ALOHA (referred to as CPODA system operation). Combinations of fixed and random access in the control subframe are also possible. The boundary between the control and information subframes is not fixed. This allows for varying the size of the control subframe according to scheduling requirements. In this way more of the communication channel can be utilized for making reservations when there are no reserved transmissions outstanding. This adaptive reservation capability of CPODA operates in the same manner as in Robert's FIFO scheme.

The number of incorporated design features allows the protocol it to satisfy intended general purpose requirements. Provision is also made for stream traffic with the reservation for stream traffic made only once, at the beginning of the stream use, and maintained by all stations in a separate stream queue. Scheduling queues are maintained by stations performing distributed scheduling. Centralized assignment can also be used when delay is not critically important. Provision also exists for priority packet scheduling.

Performance measurements conducted by Chu and Naylor [CN 78] show that in the CPODA scheme there is a tradeoff that exists between the maximum throughput achievable and the corresponding average packet delay experienced for throughput levels below the maximum. Given a minimum control subframe size, this trade off is determined by the PODA frame size. Maximum achievable throughput increases with frame size but causes longer delays since a larger frame increases the time waiting for a frame boundary to send data or control information. With the piggybacking of reservations, packet delay decreases slightly as traffic load increases before rising asymptotically at maximum throughput. The general flexibility of the PODA scheme allows it to offer good overall performance characteristics.

3.225)*Dynamic Fixed Reservation Access-Control (DFRAC) [Rub 79]*

This access schemes employs a fixed time frame structure in which each frame is divided into two periods: a reservation period, and a service period during which reserved packets are transmitted. Access to the reservation channel can be either contention-free or on a random access basis. The distinguishing feature of this scheme is its dynamic element in which the number of reservation and service slots are continually variable. The system incorporates an estimator to measure traffic flow rate (and any other relevant traffic-message statistics) and this provides the necessary information for determining the proportion of reservation and service slots.

The protocol operation is otherwise typically that of reservation access with the potential for relatively high throughput. With its adaptive mechanism this scheme can deliver good throughput-delay performance (essentially low message response times) over a wide range of traffic conditions. It however requires a perfect estimate of the underlying network traffic intensity which as a practical matter is not easily achieved under conditions of short-term fluctuations. The performance is shown in the evaluation conducted by Wieselthier and Ephremides [WE 80] (see figure in section 3.226 below).

3.226) *Asynchronous-Reservation Demand-Assignment (ARDA II) [Rub 79]*

This is a reservation access scheme designed to automatically adapt to network traffic values. As is typical of reservation schemes, channel time is divided into reservation and service periods. Reservations, however, are achieved by declaring a slot to be a reservation slot in an asynchronous manner, rather than in a fixed periodic fashion. At any time, the first slot not allocated to serve a packet is declared to be a reservation slot. The difference between this scheme and DFRAC described above, is the abscence of a fixed frame structure and a different mechanism for determining the frequency of occurence of reservation slots. Reservations are made for all packets at a station at the time of a reservation slot. Message packets are then transmitted one round trip propagation delay following their reservation in slots assigned according to the system service algorithm. It is assumed that all network terminals can make their reservations within a single contention-free slot.

One round trip time,$R$ slots, after a reservation slot (or after the last in a group of reservation slots), the remaining number packets to be transmitted in the present channel service period is observed. Depending on the number of outstanding reserved packets to be transmitted, one of two methods is used to determine the position of the next reservation slot: i) If the remaining service period is less than $R$ slots, the group of $R+1$ slots following the present service period are declared to be reservation slots. ii) If the remaining service period is greater than $R$ slots (as might be the case under high throughput conditions), the slot occurring $R$ slots prior the end of the service period is declared to be a reservation slot. Performance of ARDA II, together with that of DFRAC, is shown below.

Figure 3.3226   Delay-throughput comparison of ARDA II and DFRAC. Source [WE 80] p. 878.

3.227) *Distributed Reservation Control (DRC) protocol [GE 81]*

In this protocol, bandwidth is partitioned into a data channel and a narrow-band (frequency divided) reservation subchannel. This reservation subchannel bandwidth can be made negligible with respect to the data channel bandwidth. The essence of the DRC algorithms is that a surrogate contention process is created at the reservation subchannel level thereby permitting more efficient utilization of the wideband data channel (not unlike any other reservation approach). The scheme employs a frame format in which the data channel is slotted for the transmission of fixed length packets. The transmission frame consists of $M$ packet slots with the frame duration being greater than the maximum roundtrip propagation time. For reservation control purposes, time is sectioned by subdividing each packet slot into $L$ equally spaced segments. All reservation signals transmitted on the reservation subchannel are identical (such as a fixed frequency tone) with time of transmission being the criteria used in establishing reservation requests.

Upon packet acceptance by a node[7], a uniformly distributed random integer variable $X$ ($1 < X < ML$), is generated and the reservation signal will be generated by the node over the reservation channel during the next occurrence of the $X$th segment. All nodes which transmitted a reservation signal on any segment of a given frame will monitor the following frame to determine (by appropriate detection method) if their reservation was granted. If the reservation request is rejected (due to collision or based on the criteria of the protocol scheme), the node will reinitiate the reservation process. Each node which obtained a packet slot assignment will transmit on its assigned slot.

DRC-I

On the frame following the reservation signal transmission, the transmitting node will monitor the reservation subchannel and counts the number ($Nx$) of the segments on which the reservation signal was detected from the beginning of the frame up to the $X$th segment. If $Nx$ is less than or equal to $M$ (the number of slots per frame), on the next frame, the packet will be transmitted in the $Nx$ slot; if $Nx > M$, the reservation request is rejected.

DRC-II

This version of the protocol relies on the ability to detect, not just collision, but the number of collided signals in a reservation segment. All nodes which transmitted a reservation signal in a given frame will monitor the reservation subchannel on the following frame. Counts of singly and plurally occupied segments will be made. Singly occupied segments of counts $< M$ are allowed reserved transmission in the following frame in the corresponding slot positions. If the number of singly occupied segments are less than the number of available slots in the frame, $M$, the remaining slots are available to be used by stations whose reservations occurred in the plurally occupied segments. The determination of selected packets is again done using a random number computation. In order to increase the rate of success in the contention for available slots, the number of contending reservation request signals is constrained to be

---

[7] The described protocol [GE 81] includes an acceptance phase which is used as a means of flow control in the system. Packets are accepted at a node depending on the overall traffic rate. The exclusion of this acceptance phase does not affect the discussion of the multiaccess operation of the scheme.

no greater than the number of available slots (here the need to determine the number of collided signals becomes important). The selection of contenders is made using a calculation sequence performed at each node. All other plurally occupied segments are rejected.

Simulations show [GE 81] that certain versions of the protocol can support throughput rates in excess of 97 percent, and together with its flow control mechanism, maintain stability even under overload conditions. Access delays are tolerable (ranging from 0.5s under light load to 0.88s for saturated traffic conditions). The unique quality of this scheme lies in its implementation of the reservation channel using frequency division instead of the time division approach of the other methods examined. Though the achievable performance of the scheme is good, the practical problems of implementing its reservation detecting mechanism may outweigh the potential gains.

In concluding this section on reservation protocol schemes we can reiterate our findings. Robert's FIFO with ALOHA reservation (modified with a fixed frame [Tas 84] section 3.221) and Rubin's DFRAC (section 3.225) offer the best performance of the reservation approaches presented. The average delay is greater than that of ALOHA for very low throughput values but they possess the high throughput capability of reservation protocols. Also worthy of mention is CPODA by Jacobs *et al* (section 3.224) that has a number of general purpose features and good, flexible delay-throughput performance. This is not surprising given that CPODA is an outgrowth of FIFO. The additional operational features are facilitated at the expense of increased complexity of implementation.

As we move on to the next two sections which describe a number of mixed and adaptive protocol schemes, we will see, quite interestingly, how the basic random access and reservation approaches whose various forms have been seen above, have been moulded into protocol methods that exploit their particular characteristics.

## 3.3 Mixed Protocol Schemes

In chapter 2, we defined mixed schemes are those which combine some variations of the basic random access and reservation approaches. They achieve an enhanced performance by integrating the attributes of both methods. This generally leads to a more balanced capability across a range of network traffic conditions.

3.301) *Group Random Access (GRA) [Rub 77b]*

This protocol scheme consists of using certain channel time-periods to allow different subsets of system users to transmit their packets on a random access basis. Under the GRA discipline, a group of network terminals are provided with a periodic sequence of channel access periods, during which this group uses a random access discipline to gain access into the channel. Packets experiencing collision during a certain period will retransmit during their next access period. Different channel time periods are thus randomly access by individual groups of stations. The idea is simply a periodic time-division assignment of channel use to particular groups of users. The approach is designed such that channel time which is not assigned to specific groups, can be allocated amongst other network users again using GRA procedures or using other access control techniques.

The GRA protocol approach yields a delay-throughput performance characteristic comparable to those attained by a regular slotted ALOHA random access procedure. It also provides the flexibility of granting access to different classes of information and protocol messages.

3.302) *Integrated Random-Access Reservation multiple access (IRAR) [Rub 77a]*

Under the IRAR scheme newly arrived packets can be designated for reserved or random access. The protocol rather than having fixed periodic reservation slots, assigns the first unreserved slot following the instant at which a collision is recognized (i.e the round trip propagation delay following any collision) as a reservation slot. Access for reservations may be made either in a contention free or

random access manner. Any slot that is not reserved is declared a random access (RA) slot. In the event a collision occurs, each colliding packet is assigned for transmission by reservation. Newly arrived packets are restricted to transmitting only in random access slots.

While in the basic IRAR scheme a packet is not allowed to try more than one random access transmission, it is also possible to employ a variation of this scheme which allows for multiple random access transmission attempts. Improved performance at medium and high network traffic values are achieved by the following scheme variations:

IRAR II

If a newly arriving packet recognizes a reservation slot prior to the first available RA slot, it makes a reservation in the slot rather than attempting a later random access transmission.

IRAR III

In addition to declaring the first slot following the recognition of a collision as a reservation slot, the first unreserved slot following a service slot is also to be declared a reservation slot. This eliminates the high possibility of collision that may occur following a long service period. (Other variations on the basic scheme have been considered which simplify the analysis though leading to similar performance results).

For channels with high propagation delays, the IRAR scheme yields delay-throughput performance characteristics superior to those obtained under pure reservation and random access schemes, for medium and low network throughput values. Stability is not an issue of concern as random access contention is limited. This scheme by Rubin [Rub 77a] first demonstrated the usefulness of a more integral combining of the advantages of random access and reservation approaches. This concept represents an important step in the development of protocol schemes and is employed, in a number of different ways, in other later proposed access protocols.

Rubin's [Rub 77a] results presented below are based on the consideration of single packet messages. Comparative analysis is also provided by Wieselthier and Ephremides [WE 80]. These results indicate that for all values of throughput, the delay performance of IRAR III is as good as or better than even slotted ALOHA as well as capable of much higher throughput. The DFRAC (section 3.225) and ARDA (section 3.226) reservation schemes are also shown.



Figure 3.302a  Delay-throughput comparison of Slotted ALOHA, ARDA I, DFRAC, and IRAR I, II and III. Source [Rub 77a].



Figure 3.302b  Delay-throughput comparison of ARDA II, DFARC, and IRAR III. Source [WE 80] p. 879.

3.303) *Distributed Tree Retransmission Algorithms  [Cap 79]*

In schemes where contention and reservation are employed, retransmissions after a conflict are scheduled by random access or conflict free reservations or by the application of an appropriate algorithm. The tree algorithm is based on the observation that a contention among several active sources is completely resolved if and only if all the sources are somehow subdivided into groups such that each group contains at most one active source. The algorithm relies on a broadcast capability of the channel. While the original tree protocol (and Gallager's window protocol) assumed a zero propagation delay, it has been shown by Liu and Towsley [LT 83] that basic multiplexing approaches in interleaving the algorithms could be used to adapt the approach to the long propagation delay of satellite channels. This results in just slightly degraded performance throughput.

54

Contention resolving tree algorithms are inherently stable. They have a maximum throughput of 0.43 packets/slot and have good delay properties. It is also shown [Cap 79] that, under heavy traffic, the optimally controlled tree algorithm adaptively changes to the conventional time-division multiple access protocol. Depending on the message arrival and length statistics, reservation schemes are capable of throughput approaching one. Since, however, the tree algorithm is a collision resolution mechanism, its performance should be compared to that of the reservation channel as this is where the algorithm would be used. In such cases it offers stability over random access reservations and delay better than, or at least as good as, contention free reservation.

3.304) *Interleaved Frame Flush-Out  (IFFO) Protocols [WE 80]*

· This protocol employs a frame format where the minimum number of slots in a frame is equivalent to the maximum round trip propagation time of the satellite channel. This ensures that the reservation information generated at the beginning of a frame is received before the start of the next frame. Control is distributed, with the broadcast nature of the channel allowing each station to maintain necessary information on reservations.

The first slot of each frame, the status slot, is divided into a number of minislots equivalent to the number of stations in the network. These reservation minislots are assigned in a (contention-free) TDMA fashion. The status slot is followed by the reserved transmission slots. Since the frame has a minimum length (slots) equivalent to the round trip propagation time, any unreserved slots are used on a contention basis. If, however, the number of reserved slots exceed the minimum frame length, no new slots are added and the system operates completely on a reservation basis. Each of the IFFO protocols is characterized by a different transmission procedure in the contention slots.

*Pure Reservation IFFO (PR-IFFO)*

Reservations are made for packets arriving at individual stations. In the event that there are unreserved slots within a frame (given the minimum frame length requirement), they are not used for

contention but simply remain idle. The operation is thus the same as Robert's FIFO reservation scheme (section 3.221) with TDMA reservations.

*Fixed Contention IFFO (F-IFFO)*

A packet arriving at a station during the interval of a contention (unreserved) slot is transmitted in the following slot. Packets arriving during the interval of reserved slots are not allowed to contend but instead are queued at the station and a reservation is made for transmission (in the second frame following the frame of arrival). Packets arriving at the last slot of a frame are not transmitted in the following slot (which is the status slot), but wait and make reservations for later transmission. This fixed contention operation is similar in concept to that employed in the RR reservation scheme (section 3.222 above) where unused TDMA assigned slots are contended for.

· *Controlled Contention IFFO (CR-IFFO)*

In this protocol version, the transmission procedure is state-dependent. In each contention slot, a station with a packet to transmit does so on the basis of a probabilistic formulation. The probability is a function of: a) the number of packets for transmission at the particular station, b) the number of the slot being contended for, and, c) the number of slots reserved in the particular frame. The transmission procedure is thus a randomized one. A heuristic policy is derived to determine the probabilistic model for transmission control.

The performance of the IFFO protocols are shown below together with a comparison of the schemes reviewed in the previous two sections. From this analysis carried out by Wieselthier and Ephremides [WE 80] we can see that the performance of F-IFFO is superior to the other schemes considered. Since PR-IFFO is in fact the same as that of Roberts FIFO with TDMA reservation, we can also use the performance of PR-IFFO to relate the performances of these mixed schemes with the reservation schemes compared earlier in section 3.221 [Tas 84]. From this comparison, we can conclude that of the schemes examined so far, IRAR III and F-IFFO offer the best overall performance.

**Figure 3.304a** Delay-throughput comparison of PR-IFFO, F-IFFO and CR-IFFO (simulated). Source [WE 80] p. 878.

**Figure 3.304b** Delay-throughput comparison of Slotted ALOHA, DFRAC, ARDA II IRAR III and F-IFFO. Source [WE 80] p. 878.

3.305) *Dynamic Reservation Multiple Access (RMA) [BP 79, GSS 82, TC 86]*

This scheme is a fixed TDMA based approach upon which a reservation mechanism is imposed. Time is divided into frames, each frame consisting of reservation slots, preassigned TDMA slots for each of the stations in the network, and access slots for the transmission of reserved packets. The reservation slots at the beginning of each frame are divided into minislots equivalent to the number of connected stations. Reservation access is thus contention free. Frame length is variable, with the number of reserved access slots at the end of each frame being determined by the number of packets reserved (in the previous frame) for transmission.

At the beginning of the frame, each station makes a determination of the number of packets it has awaiting transmission. One is subtracted from this number (since it will be sent in the station's preassigned TDMA slot), and a reservation is made for the remaining number of packets using a transmitted reservation vector. Immediately after the stations have made their reservations, each transmits in its preassigned slot. The total length of the preassigned section of the frame is greater than the roundtrip propagation delay. As a result, the reservation requests of all stations are conveyed

57

across the network by the end of the preassigned slots. The asymmetry of the preassigned TDMA slots with regards to terminals means that the transmission time for packets at terminal $m$ always occurs after that of station $m-1$. By constantly rotating the order of the slot allocations, this problem is however overcome.

The operation of this scheme is such that it produces a performance significantly better than that achieved by the ALOHA method and earlier schemes by Binder [Bin 75] (sect. 3.222) and even Roberts FIFO scheme [Rob 73] (sect. 3.221). Guha *et al* [GSS 82], obtained improved performance by using an algorithm that incorporated prior reservation requests, in addition to received reservation vectors, in granting reserved slot assignments. This was also further refined by Tsai and Chang [TC 86], again with a different handling of channel reservation in which a greater number of reservation request bits would be used, or the two bits (used by Guha) maximally utilized, in determining reserved assignments. This improved performance was more pronounced at higher throughput values (see below - the vertical axis in Fig. 3 [TC 86] is a measure of Q, the average number of packets waiting in queue at a station at the beginning of a frame).



Figure 3.305a    Delay-throughput comparison of ALOHA, Roberts R-ALOHA, Binder's RR, TDMA, and RMA. Source [GSS 82] p. 59.

Figure 3.305b    Comparison of the number of queued messages versus throughput for RMA scheme and proposed improvement by Tsai and Chang. Source [TC 86] p. 727.

3.306) *Combined Random/Reservation Multiple Access (CRRMA) [LM 83]*

This method uses an on-board satellite processor which provides centralized coordination of the scheme. Channel time is slotted with each fixed length slot containing a data portion as well as a reservation segment. The data portion of the slot operates in one of two states: a Contention state, or a Reservation state. The transmission slot header (reservation segment) of the up-link channel is divided into $N$ minislots each of which is long enough to accommodate a request packet. The number of minislots, however, are an order of magnitude less that the number of stations in the network[8]. Collisions may thus occur between request packet transmissions. In such an event, the slot is then in a retransmission (RETX) state. If request packets in a given slot are successful, the slot is in a FREE state. For each data packet at a station, a request packet is transmitted in one of the $N$ minislots chosen at random. Transmission of request packets and corresponding data packets may take place independently depending on the state of the transmission slot (Request-Data segments). The possible states are RETX-Contention, RETX-Reserved, FREE-Contention and FREE-Reserved.

In the RETX-Contention state, it is likely that a large number of request packets have collided. The data portion of the slot is thus unusable for Contention retransmissions since there would be certainty of collision. The data portion in such a case is partitioned into $N'$ additional request minislots. This expedient assists in reducing the backlog of collided reservation request packets. The scheduler examines the successfully received requests as well as the data slot transmissions and makes request assignments. The outcome of any transmission, request and/or data packet, is known at the user stations after the round trip propagation time of $Q$ slots. Two versions of the protocol are proposed:

---

[8]    While this introduces the possibility of instability in the operation of the scheme, the overhead associated with packet reservation is not proportional to the user population size. This issue of instability is addressed by the Controlled Channel Access (CCA) versions of the scheme.

*Uncontrolled Channel Access (UCA)*

A station transmits its request packet in one of the N minislots chosen at random. If the data slot portion is in the contention state, it transmits a copy of the tagged packet. The results of these transmissions are: i) data transmission is successful - i.e. data slot was in the contention state and there was only one request packet in the minislots, ii) the request was successful, and iii) the request collided. Based on these results, the actions taken Q slots later are: i) Data transmission successful: discard copy of tagged packet at transmitting station. ii) Request successful: wait and transmit packet in assigned slot. iii) Request collided: retransmit request packet in one of the N minislots of next slot.

*Controlled Channel Access (CCA)*

The first request for a packet at a station is made only in a FREE slot. New packet requests are not transmitted following a request collision. This reduces the chance of further request collisions. Request packet collisions are resolved according to a predefined algorithm. This eliminates the possibility of instability on the reservation channel. Transmission in a given slot depends on the state of the slot. i) RETX-Contention: transmit a packet in one of the N (or $N + N'$) minislots. ii) FREE-Contention: transmit a request packet in one of N minislots and at the same time transmit a copy of the tagged packet in the data portion. iii) Reserved (RETX or FREE): transmit a request packet in one of the N minislots. (The station assigned the data portion transmits in the data slot).


The UCA scheme has an inherent (random access) stability problem which becomes more prevalent as the number of minislots N decreases. This problem is however addressed by the CCA scheme. As seen in the figures below, the CRRMA protocols exhibits good delay-throughput characteristics over a wide range of utilization. For the number of users greater than 5, both schemes give similar performance (UCA being simpler however). The use of an on-board scheduler reduces the round trip delay associated with receiving reservation assignments. Like the other integrated reservation/random access approaches above, the performance of this scheme is comparable to that of slotted ALOHA at low throughput yet able to achieve much higher maximum throughput values. ($\theta$ in the figures below, specifies the faction of reservation overhead).

**Figure 3.306a**  Delay-throughput comparison of CRRMA-CCA with Slotted and Reservation ALOHA. Source [LM 83]. p. 1170



**Figure 3.306b**  Delay-throughput comparison of CRRMA-UCA with Slotted ALOHA and SRMA. Source [LM 83] p. 1170.

**3.307)** *Announced Retransmission Random Access (ARRA) Protocols [Ray 85]*

This contention based scheme increases the capacity of a slotted random access channel by adding a small amount of control information to each message transmission so that future collisions involving retransmitted packets can be avoided. The only source of contention in this scheme is between packets transmitted for the first time. The transmission slots in this protocol are grouped into frames of $K$ slots. Each slot is divided into a message slot portion and a minislot portion consisting of $K$ minislots. At the start of each frame is a group of $K$ minislots - the common minislot pool (CMP).

In the basic ARRA protocol, each message transmitted in a given slot is accompanied by an announcement within one of the corresponding minislots stating the retransmission slot number $(1 - K)$ to be used in the event of a collision. This slot number is randomly chosen. Users then monitor both the message and announcement channels (which includes the CMP), storing one frame's worth of channel

61

feedback. At the end of a frame, based on the outcome of the attempted transmissions and the retransmission announcements, each user is able to determine the slots that will be available for new packet transmissions in the next frame.

Packets arriving at a station during a given frame are held until the end of the frame at which point a determination can be made of available slots for transmission. The packets are then transmitted at random in any of the available slots in the next frame. If based on the collisions occurring in a particular frame, there are no available slots in the next frame for transmission of new packets, the new transmission cannot take place[9]. In such a case the new transmissions are replaced by an announcement in the appropriate slot of the CMP for transmission in a future frame. This protocol operation thus avoids all collisions between new and retransmitted packets. However, because packets are retransmitted unconditionally, conflicts may occur between packets that have announced the same slot for retransmission.

An extended ARRA protocol is also proposed which reduces predictable collisions by changing the unconditional nature of packet retransmissions. In this protocol additional constraints are provided so that collided packets that have announced the same retransmission slots are prevented from attempting to transmit. This avoids certain collision while also increasing the pool of available slots for the transmission of new packets. The aborted packets use the CMP to announce for a later transmission.

The performance characteristic of this protocol is very similar to that of slotted ALOHA while offering a greater maximum throughput and lower delays. Basic ARRA, with unconditional retransmissions offers a capacity of about 0.53 while that of the extended scheme is about 0.6. This

---

[9]     A straight forward extension allows retransmissions to be made randomly over M>1 future frames.

performance is however achieved at the cost of greater complexity as well as the practical difficulty of properly maintaining synchronization in order to isolate minislot boundaries within transmission slots.



Figure 3.307a   Average number of retransmissions vs throughput comparison for Slotted ALOHA,Basic ARRA and Extended ARRA.

Figure 3.307b   Delay-throughput comparison of Slotted ALOHA, Basic ARRA, and Extended ARRA. Source [Ray 85] p. 1188.

3.308) *Scheduled-Retransmission Multiaccess (SRM)* [10]*Protocol [Yum 87]*

This scheme is similar to the ARRA approach above and based on random access contention, with collided packets rescheduled to limit further collisions with new or other retransmitted packets. This is a centrally controlled scheme with satellite processing providing the necessary arbitration unlike the ARRA, though the operating concept is the same.

In the Fixed Frame (FF) version of this protocol, the satellite channel is divided into frames of a given number of slots. Each frame is divided into an ALOHA subframe and a reservation subframe. Each slot in the ALOHA subframe has a body and a header. The body of each slot accommodates one packet. The header is divided into minislots one of which is set at random and used for transmission rescheduling purposes in the event of a collision. For a collision in a given slot, the header is examined and the

---

[10]   We refer to this scheme as the SRM rather than its given acronym SRMA which is the same as that of the Split-Channel Reservation Multiple Access (SRMA) scheme [TK 76] proposed at an earlier date.

particular minislots set by the colliding packets used to differentiate them and order retransmissions. If two packets have set the same minislot, the information is garbled and both packets are unable to schedule a reserved retransmission. They will instead be retransmitted in a future ALOHA slot.

A new packet arriving at a contention (ALOHA) subframe is immediately transmitted. At the end of a period given by the round-trip propagation delay the outcome of the transmission as well as future reservation decisions are known. If unsuccessful, the packet is scheduled for a reserved retransmission by assignment to a dedicated slot in the reservation subframe which guarantees success. If the packet is also unsuccessful in making a reservation, due to collisions in the header minislots, it is aborted and retransmitted in a future non-reserved slot, chosen at random. Packets arriving at a reservation subframe are scheduled for retransmission at random amongst the next sequence of ALOHA subframe slots. Because of the use of fixed length frames it may also be necessary to abort packets that cannot be accommodated in a reservation subframe.

In the Dynamic Frame (DF) version of the protocol, the framing pattern is different. Here the ALOHA frame (not a subframe) is of fixed length with the reservation frame of variable length to accommodate all successfully reserved packets.

The Fixed Frame SRM can achieve a maximum throughput of 0.665 and the Dynamic Frame version at least 0.93. Both protocols give average delay considerably lower than that of slotted ALOHA even when throughput is as low as 0.2. The protocols can also incorporate a priority scheme that can provide different average delay for different priority users. The protocol is not unconditionally stable but can be made so by adjusting the retransmission delay parameter of aborted packets dynamically according to traffic conditions. This protocol suffers some of the same practical operating difficulties as the ARRA scheme above.

**Figure 3.308a** Throughput versus offeredtraffic comparison of Slotted ALOHA,Fixed Frame SRM, and Dynamic Frame SRM. Source [Yum 87] p. 8.7.5.



**Figure 3.308b** Delay-throughput comparison of Slotted ALOHA, Extended ARRA, FF-SRM and DF-SRM. Source [Yum 87] p. 8.7.5.

As we have seen from the mixed protocol schemes presented above, the proper integration of random access and reservation approaches can lead to the achievement of superior delay-throughput performances. The coordination of that integration, however, introduces an additional level of complexity. In the latest schemes examined, Raychaudhuri's ARRA (section 3.307) and Yum's SRM (section 3.308), the inclusion of control information in the random access phase, though practically difficult to achieve, leads to a very successful analytic solution. The other high performance schemes are Rubin's IRAR III (section 3.302), Wieselthier and Ephremides' F-IFFO (section 3.304), and Guha's *et al* RMA (section 3.305) which offer high throughput values and low average delays. Of these access methods, F-IFFO and IRAR III are the most easily implemented. The key performance difference between these mixed schemes and the most successful reservation approaches, such as Robert's FIFO, is their ability to achieve very small average delays at low throughput yet still offer high throughput capabilities.

## 3.4 Adaptive Protocol Schemes

Adaptive protocol schemes provide flexible operation over a range of operating conditions by dynamically adjusting to meet the needs of network users. In an environment where the demand for the communication channel is bursty, a scheme that is not adaptive may result in severe waste of the channel capacity. The schemes detailed below possess the ability to change their method of access, to some extent, with varying load and traffic situations. Essentially, adapting from the extremes of random access through reservation access and eventually to TDMA operation, as demanded by network loading increases. These protocols are thus able to changes their operating mode to deliver the type of access operation that is most efficiently suited to the network conditions encountered. These adaptive schemes are therefore able to provide high performance over a wide range of operating conditions.

### 3.401) *URN Scheme [KY 78]*

This scheme is proposed in the context of a fully connected broadcast environment. The time axis is divided into packet slots and users are synchronized. The scheme works by giving the full right of access to the channel (i.e right to transmit with probability 1) to $k$ out of the $M$ users in the system. A successful transmission occurs if there is only one busy user (i.e a user with packets to transmit) in the group of $k$ selected. The probability of this successful transmission event is maximized for $k = [M/n]$ (the integer value of), where n is the number of busy users. In a lightly loaded system with n=1, $k=M$, and so all users are allowed to transmit. Hopefully only one will do so, the one busy station. At the other extreme, if the system is fully loaded with n=$M$, then $k=1$ and only one user is allowed to transmit. The scheme thus converges to TDMA. It is assumed that all users have knowledge of the number of busy users, n. The necessary sampling of $k$ could be done using two different approaches, either in a random manner or without repetition from slot to slot until all users are sampled once. This would correspond to a convergence to either random TDMA or round-robin TDMA.

This protocol relies on estimating n, and selecting the *k* users who are given access rights to the channel in a given time slot. One means of determining $n$ is the use of a single reservation minislot at the beginning of each data slot which can be used by stations to indicate their going busy. This is done in conjunction with a monitoring of stations going idle through packet acknowledgements. This scheme, however, only allows for the detection of 0, 1 or >1 (collision on minislot) reservation. There is thus estimation errors in determining $n>1$. However, it is shown through measurement and simulation that $n$ can be closely approximated using the single minislot mechanism. In addition, the protocol is shown to be insensitive to small errors in the determination of $n$ by users. Other heuristic algorithms for estimating $n$ have also been studied by Lam and Kleinrock [LK 75].

Selection of the group of *k* users can be achieved with pseudorandom generators. These generators at all stations allow them to draw the same pseudorandom numbers. Another approach, referred to as the round-robin slot sharing window mechanism, consists of having a window of size *k* move over the population space. If a collision occurs, the window downsizes. If no collision occurs, the window is advanced to cover the next set of previously uncovered users; the size of *k* being determined by $n$. Such windowing schemes adapt to both the total load and the results of the selections of the previous slot. The selection of users for access rights in a given slot could also employ some priority mechanism.

The operation of this scheme is such that it smoothly adapts to the load on the channel, varying from slotted ALOHA at light load, ranging through an asymmetric scheduling scheme, finally eliminating collisions as it converges to TDMA in the heavy traffic case. The channel control overhead to achieve this operation is negligible and does not depend on system size. The important need, however, for traffic measurement information and distributed coordination in the group selection of users, presents difficult issues in the implementation of this scheme. The adaptive performance of the scheme is shown in the figure below.

**Figure 3.401**    Delay-throughput comparison of Optimal ALOHA, TDMA, URN, and Perfect scheduling. Source [Tob 80] p. 484.

3.402) *Split Reservation Upon Collision (SRUC) [BF 78, TI 86]*

This protocol consists of a contention scheme with a superimposed reservation technique which automatically operates when collisions occur. In this Reservation Upon Collision (RUC) scheme, channel time is divided into fixed length slots with these slots grouped into fixed length frames[11] . Each slot consists of a data packet and header portion together with a subslot portion divided up into a number of minislots and used for the conveying of signaling (control) information. These signaling information (SI) minislots are assigned to stations in a TDMA (or other orthogonal multiplexing) fashion. The number of SI minislots per transmission slot is determined such that over an entire frame each station has an assigned minislot.

The data subchannel exists in one of two states: the contention state or the reservation state. It is normally in the contention state. In this state, each user can continuously transmit packets under a slotted ALOHA protocol as long as no collision occurs. In addition to transmitting packets, each station

---

11    To achieve effective performance with time-varying traffic, it was shown [TI 86] to be necessary to select an optimum frame length according to network loading. For the purpose of implementing the adaptive selection of the frame length, the V-SRUC (variable frame length SRUC) protocol was proposed in [IT 82].

transmits in its own SI minislot information on the outstanding packets which the user has to the time of the transmission of it SI. This information pertains to untransmitted packets at the station as well as packets for which feedback confirmation of its transmission outcome is not yet known.

When a collision occurs the data subchannel switches to the reserved state. In this state, a Channel Control Procedure is implemented which uses the information available in the SI minislots to set up a common reservation queue for the transmission of packets (header information is used to monitor the successfully transmitted packets and remove them from future consideration). The system then continues in this reserved state until all backlogged packets are cleared. In this protocol therefore, newly arriving packets at a station are immediately transmitted if upon arrival the data subchannel is in a Contention state, otherwise it is delayed, as are the collided packets, until it is assigned a reserved slot for transmission.

From the ability of the data subchannel to operate in a reserved state, the protocol is stable under all traffic conditions. Performance analysis shows the advantages of the SRUC technique are mainly due to its split capability which reduces the degradation of the performances when the number of users increases. Because reservation information is sent as part of each frame, a roundtrip propagation delay is avoided in the reserved retransmission of collided packets. Further analysis conducted by Tasaka and Ishida [TI 84], for SRUC with multiple-packet messages and go-back-n (GBN) ARQ, have also confirmed the utility of the scheme (see figure below).

Comparing the results here with that of the URN scheme in the section above, it is seen that URN is able to provide better delay performance for low throughput values, better emulating that of slotted ALOHA. The SRUC protocol, however, does not offer as much an implementation problem as the URN scheme.

**Figure 3.402a** Delay-throughput comparison of Slotted ALOHA and SRUC with go-back-N ARQ. Source [TI 86] p. 943.

**Figure 3.402b** Delay-throughput comparison of S-ALOHA (GBN) Res ALOHA, Pure Res. with ALOHA res., Pure Res. with TDMA res., & SRUC. Source [Tas 84]. p1578.

3.403) *Random Access with Notification (RAN) [CBK 89]*

This scheme follows in the trend of overlaying random access operation with the scheduling of retransmissions in the event of collisions. The protocol, is uniquely designed for operation in a VSAT network environment with the communication paths forming a star connectivity with no requirement for direct communications between individual stations[12]. Another notable assumption of this protocol is that stations are not able to listen to their own packet transmissions. They thus rely on the central hub for a determination of packet collisions. The hub additionally monitors traffic activity and adapts system parameters to traffic changes. A separate signaling channel is used by the individual terminals for communication with the central hub. This channel is operated on a contention-free TDMA basis.

The data channel consists of two types of slots: random-access (RA) and reservation (RES), with collided and erred packets prevented from transmitting in the random access slots. New packets use the

---

[12] This configuration allows for the design of inexpensive VSAT terminals that communicate with a larger hub terminal.

RA slots for initial transmission. The hub monitors the traffic rate generated by the stations and for a given delay response requirement, determines the key system parameters of the RAN protocol, namely the mix of RA and RES slots. When traffic exceeds a certain intensity the hub issues a command for terminals to switch over to full reservation operation. This is switched back after the traffic load has fallen back to acceptable levels.

For a packet generated by a terminal, if the next data slot is an RA slot, the packet is transmitted in that slot. If the next slot is reserved, the packet is randomly scheduled for transmission among the next available RA slots (excluding those that have already been assigned a packet by that terminal). The station keeps a copy of the packet along with the frame number and the time slot used for transmission. This information allows for retransmission of the particular packet after determination of collision or error by the hub. Over the signaling channel terminals notify the hub of the number of packets transmitted. The hub does the necessary processing and for each terminal it correlates the information received on the signaling channel with its observation of packets on the data channel. If the number of correctly received packet are less than that notified by a station, the hub allots the required number of slots to the terminal for contention-free transmission in a future frame. The RES slots are randomly ordered and scheduled among the stations in such a manner that the stations receive notification in time to retransmit in the allocated reservation slot. There is continual return acknowledgement is this RAN protocol.

When the number of packets lost (collided or erred) through random access transmission in a particular frame exceeds the number of RES slots per frame, for a certain number of consecutive frame, the hub broadcasts a message for terminals to switch over to reservation TDMA, starting with a specific frame number. The hub then completely controls the scheduling of transmissions. When the number of packet reservations per frame falls below the threshold number, for a consecutive number of frames, a command is issued for a change back to random access operation.

When not operating in the reservation mode, the delay-throughput performance of this protocol is typical of contention based schemes where delay is low under conditions of low traffic load, increasing as throughput increases. The RAN protocol is able to provide higher throughput than slotted ALOHA, however, delay may be greater given the necessary signaling overhead. This is dependent on the number of terminals in the network and the amount of signaling required. When throughput is low the delay characteristics far outperform those of reservation schemes because the packets do not have to go through the negotiation process for the allocation of transmission. As throughput increases, the delay begins to approach those of reservation schemes. The key advantage of RAN is its ability to adapt across a range of operating conditions though a great amount of coordination overhead is required.

## 3.5 General Conclusions from Protocol Examination

At the end of chapter 2 we sought to establish a basic framework for the judging of performance amongst different protocol schemes. In spite of the many other important attributes, the main focus of our study was to look at the channel throughput versus average message delay tradeoff. As stated, the objective of any access protocol is to achieve a maximum channel utilization, in keeping with traffic demands, and subject to given packet delay requirements. Given the bursty nature of the communication environment, demand assignment rather than fixed assignment schemes are appropriate. Of these, contention protocols give faster response times when the channel is lightly loaded, but are inherently unstable. Even though stabilizing control mechanisms can be employed, throughput rates do not exceed a small fraction of the channel capacity. Reservation protocols on the other hand, offer inherent stability but suffer the disadvantage of slower response times even when the traffic load on the network is very low, as reservations for transmissions are still required. Given the respective performances limitations of either of these two basic classes of protocols, the key to achieving enhanced overall performance is in the appropriate combining of these approaches.

With bursty data communication, one extreme scenario is that in which no channel access is being requested. From this extreme situation it is desired that a station needing access for a newly arrived packet is provided with immediate transmission rights. This ensures that a minimum transmission delay is achieved; the issue of channel utilization not being a factor given the very low throughput demanded. With the arrival of new packets at a station being a stochastic event, achieving such an instantaneous access response requires that the protocol scheme be one of random access in which no control information is necessary prior to a transmission.

On the other operating extreme is the case of the fully loaded network in which each station has packets ready for transmission. Neglecting the issue of priority, the overriding desire is to have packets transmitted by each station without the possibility of collision. This would insure that the

channel is utilized to a maximum. The appropriate type of access control would be a pure TDM access in which users were allowed to transmit in order, and in which there is no requirement for reservations. This type of operation would also maintain an evenly distributed minimum delay.

Given these two operating extremes, a protocol scheme designed for efficiency across the spectrum of network conditions must possess the ability to adapt to the access modes demanded. The transition should also be one that appropriately addresses the intermediate traffic conditions. From the preceding examination of protocol schemes we have seen that the innovative integration of random access and reservation methods have been most successful in achieving this desired objectives.

In concluding section 3.2 on reservation schemes, we stated that within the scope of performance as measured by the delay-throughput, the best reservation access protocols were Robert's FIFO with ALOHA reservation (modified with a fixed frame [Tas 84] section 3.221), and Rubin's DFRAC (section 3.225). The difference in these being the frequency of occurence of reservation periods. In DFRAC it is continually variable wheras in FIFO the reservation slots are periodic or occupying the entire channel. Their average delays are greater than that of ALOHA for very low throughput values but they possess the high throughput capability of reservation protocols. CPODA by Jacobs *et al* (section 3.224), an outgrowth of FIFO (with its general purpose features and flexible delay-throughput characteristic), also provides good performance.

These reservation schemes, though involving some amount of random access, were eclipsed by the mixed schemes that took a more integrated random access/reservation approach. The results of section 3.3 have shown the success of schemes such as Rubin's IRAR III (section 3.302), Wieselthier and Ephremides' F-IFFO (section 3.304), and Guha *et al*, RMA (section 3.305) protocol, as well as later protocols by Raychudari, ARRA (section 3.306), and Yum, SRM (section 3.307). The latter two not only involved the integration of random access and reservations but did so by introducing control information into the random access communication. By so doing the time taken to obtain a scheduled retransmission

in the event of collision was effectively reduced. These schemes, however, required greater synchronization efficiencies and introduced further complexity in the operation of the protocol. As a result, IRAR III and F-IFFO can be considered superior, given their relatively more simple means of implementation. Of the adaptive methods, both Kleinrock and Yemini's URN (section 3.401) and Borgonovo and Fratta's SRUC (section 3.402) were able to deliver good performances over the total throughput range.

The key concept of mixed and adaptive schemes based on an integrated approach to employing reservation when random access becomes less ideal, is one clearly demonstrated by this study. The effectiveness of this concept has also been further extended through the conveying of control information during random access contention to enhance the reservation response times in the event of collisions (ARRA and SRM).

Other useful observations that have been made include the fact that: i) random access operation though limited as a stand alone means of operation has the attribute of fairness as well as simplicity and ease of implementation. ii) reservation schemes, either controlled centrally or distributedly, possess the capacity for prioritizing transmissions. This however requires an additional information overhead which further reduces the channel capacity devoted to actual data transmissions. iii) the majority of schemes proposed employ distributed control which additionally demands the accurate maintaining of distributed global queues. iv) implicit reservation schemes or protocols that incorporate this type of reservation are most effective for multiple packet messages and stream traffic, though they may suffer from fairness deficiencies. v) fixed structures may allow optimization of a protocol for certain throughput ranges but unless there is flexibility in the assignment of reservation and service segments the scheme is limited in its ability to deliver good performance over a wide operating range.

Even within the limited scope of the throughput-delay measure, it should be apparent that our assessing of schemes could only be done in very general terms. From the examination that preceded it is

difficult to single out one scheme as being optimally superior. Nonetheless the following complexity-performance matrix can be defined for the more successful schemes:



IRAR - IRAR III (section 3.302)
IFFO - F-IFFO (section 3.304)
SRUC (section 3.402)
URN (section 3.401)
RMA (section 3.305)
ARRA (section 3.307)

**Figure 3.5.1**     Performance-Complexity tradeoff comparison for the most successful protocol schemes examined.

In addition to general delay-throughput response, the performance given above is defined in terms of minimum average delay as well as maximum throughput achievable.

In the next chapter we examine some of the parallels between the satellite access schemes presented above and some of the recently proposed optical access control methods. Confirming the application of satellite protocol approaches to the optical problem, we will use the understanding gained from this section of our study to provide the basis for a new optical access protocol.

# Chapter 4

# Examination of Multiaccess Optical Network Protocols

## 4.1 Introduction to Optical protocol methods

As initially introduced in chapter 1, there are certain characteristics of the high speed optical network environment which are similar to those of the satellite networks; the key element being the large normalized propagation delay encountered. That is, as the communication capacity of frequency channels get larger in optical networks, the ratio between end-to-end propagation delay and packet transmission time accordingly increases. It was in fact this premise that provided the thrust for our approach in studying and analyzing the various schemes proposed for multiaccess communication in the satellite environment. In conventional single channel, broadcast systems, this significant propagation delay restricts the type of access control protocols that can be employed and hence the consequent utilization of system bandwidth. Any reliance on multiaccess schemes such as carrier sensing results in a rapid deterioration of overall performance as propagation delay increases.

Maximizing the communication potential of optical fibers has been based on a wavelength division multiplexing (WDM) approach. The large optical communication bandwidth is divided up into a number of parallel communication channels, each of these acting in a manner similar to that of a single satellite channel. The access methods must now, however, be two-dimensional. Not only is there

multiaccess contention on each channel but there is also multiaccess contention for the use of the particular frequency channels (which in an unrestricted network will usually be less than the number of connected users). This two-dimensionality of course presents unique coordination problems in deriving an access control protocol.

Communicating in this multichannel environment, network stations must be capable of variable frequency operation if each station is to have the flexibility to interconnect directly with any other. This raises the very important practical concern of achieving sufficiently agile transmitters or receivers that are able to tune over a wide range of frequencies while doing so in a very short period of time. To date, those receivers best able to tune over a large frequency range are limited by their slower speeds of operation. The use of single tunable transceivers at each network station represents a design minimum and is the most practical implementation since it reduces hardware duplication and complexity.

In some network configurations, stations may have fixed frequency transmitters which then forces the use of a single communication channel by each station. In such cases, this limits the number of possible users in the network to the number of available fiber communication frequencies. One dimension of multiaccess contention is thus removed, with contention only remaining at the level of station receivers (destination contention) where multiple transmitters may desire access to a single receiver[13]. While this reduces the complexity of the station transmitters, it does not however remove the need for fast tunable receivers. Given the large available communication bandwidth of optical fibres, the restriction on the number of possible network users caused by such configurations may not be overly confining and may still be in the hundreds.

---

[13]    A similar situation is encountered in the case of tunable transmitters with receivers fixed to a single frequency channel (the contention is now at the channel level).

Another class of multichannel communication approaches is that which employs a multihop topology with a fixed set of routes assigned between a source and a destination and with buffering capabilities at each hop. A deterministic fixed store-and-forward routing strategy is used over the hops. The advantage of such approaches is that it alleviates the need for tunable lasers. However, a significant fraction of electronic speed is lost to data forwarding. This multihop operation is not as useful for data communication where there may be very short delay requirements. In our study we will tend to focus on the more general, more flexible configurations utilizing tunable components.

In the next section we examine the operation of three proposals for multiaccess communication based on the optical multichannel architecture. These schemes incorporate the concepts of satellite multiaccess packet communication detailed in our earlier chapters. The presentation of these methods provides some examples of the operational concerns involved in multiaccess communication within multiple channel systems and provide a background for the discussion of design issues that follows.

## 4.2 Proposed Optical communication protocols

4.21) *Protocols for Very High-Speed Optical Networks using a Passive Star topology [HKS 87, Meh 90]]*
The schemes presented here are extracted from the work on LAN optical protocols proposed by Habbab *et al.* Only those methods which are insensitive to propagation delay are considered[14]. In these protocol schemes the optical communication bandwidth is wavelength divided and each network station is equipped with a single tunable frequency transmitter and a single tunable receiver. A common control channel (of specified wavelength) is used by the stations to notify of their intent to transmit to a

---

14    Because LANs are of small diameter, the time taken for signals to propagate across the networks is still sufficiently short relative to transmission time to allow propagation delay-sensitive protocols to find useful application.

particular receiver and to indicate the frequency channel on which the transmission is to occur. A passive star coupler is used as a broadcast medium to cross-connect users.

A station with a packet to transmit first sends a control packet on the control channel which contains the transmitter address, the receiver address and the data channel wavelength to be used. This is similar in principle to the reservation operation encountered in satellite access schemes. The control packets are of length much smaller than that of the data packets to be sent. All idle receivers tune to the common control channel and listen for their address. Upon recognizing it, they will immediately tune to the particular frequency indicated. Acknowledgements of success and failure of control and data packets is obtained by each station listening to the echo of own packets which is facilitated by the broadcast nature of the star coupler. If a collision occurs on either the control or data channel, the user waits some randomized delay before retransmitting the original control and data packet.

Versions of the protocol are distinguished by the combinations of access methods used on the control channel and the data transmission channels.

*ALOHA/ALOHA*

Each stations transmits a control packet over the control channel and immediately follows with its data packet on one of the data channels chosen at random. The throughput on either channel is that of unslotted ALOHA. If a collision of the control packet occurs a retransmission follows after some randomized delay.

*Slotted  ALOHA/ALOHA*

Each station with a packet to transmit transmits a control packet at the beginning of a preassigned slot time and follows with a data packet on one of the data channels chosen at random in the next slot time.

*Slotted  ALOHA/N-Server  Switch*

Each station with a packet to transmit monitors the control channel over a period of one data packet length. From this, each transmitter can determine exactly which data channels and which receivers

are idle. A control packet is then sent on the control channel (according to slotted ALOHA) identifying the data channel to be used. If all channels are seen to be busy, the packet is blocked.

*Improved Protocol Operation*

Mehravari [Meh 90] provided a means of improving the throughput performance of the above protocol by introducing a simple additional requirement to the operation. He recognized that with the given operation, a collision occurring on the control channel meant that the immediate transmission of a data packet would be of no use since notification would not be available to the intended recipient. Bandwidth is thus wasted when the data packet is transmitted on the data channel if the corresponding control packet is lost on the control channel. The improvement to the scheme was to make the data packet transmission contingent on the success of the preceding control packet. The protocol operation would thus be to transmit immediately on one of the data channels at random *if* the control packet was successful.

Not acknowledged in this scheme, however, is the fact that waiting for the outcome of the control packet transmission introduces a round trip propagation delay in the transmission of the data packet. Since the environment is one of significant propagation delay, this added waiting time greatly affects the overall average delay experienced by data packets. System throughput is thus increased, but only at the expense of larger delays. In addition, introducing such delays into a random access scheme could lead to very unacceptable performance.

4.22) *Frequency-Time Division Multichannel Allocation (FTDMA) Protocols [CG 88b]*
This protocol is also based on a multichannel architecture derived through frequency division. Channel access is governed by a fixed cycle whereby transmission rights are allocated in each slot to source-destination pairs on one of the frequency subchannels.

The network model is one in which there are $N$ nodes and each of these nodes possess $N-1$ buffers that store packets for transmission to each of the other nodes, and 1 buffer for receiving packets. Thus $N-1$ virtual users are considered to exist at each node. The system can be viewed as an $N \times N$ matrix with elements corresponding to possible source-destination communication links. In each cycle, source-destination oriented permissions are granted specifying for each node or virtual user, the channel and the slots in which to transmit. Each node interface is implemented with a single receive/transmit mechanism. The class of protocols differ according to the number of permission granted in a slot and the channel allocation policy - random or deterministic. It is assumed that the choice of packets for permission assignments is made prior to each slot.

*Source/Destination Allocation Protocol*

This protocol version avoids collisions and destination conflicts by assigning permissions to disjoint source/destination pairs equal only to the number of frequency channels, $b$, in the network. Frequency channel and source/destination is assigned. A single cycle is complete when each source/destination pair (matrix element) has been allowed an assigned transmission. This is a fixed assignment approach therefore lacks the potential efficiency of demand assignment. Performance is reflective of the TDMA type operation.

*Destination Allocation Protocol*

In this protocol the number of permissions per transmission slot is increased granting permissions to $N$ users under the condition that no destination conflicts occur. In the $(N \times N)$source/destination matrix this corresponds to choosing N elements such that no column of the matrix has more than one selected element. Channels are chosen at random and since $b < N$, it is possible for channel collisions to occur in a given transmission slot. Performance is characteristic of random access schemes with throughput rising to about 0.25 for $b = N$.

*Source Allocation Protocol*

Potential access is further increased by granting $b$ permissions to source nodes for each to transmit on a different channel. In this case the number of users (real and virtual) that can potentially transmit is

equal to $b(N-1)$. This corresponds to selecting $b$ rows of the source/destination matrix and allowing transmissions by any of the $b$ sources selected. Even though transmission permissions have been increased beyond that provided in the Destination Allocation protocol, only destination collisions can occur. Performance characteristic is similar to that of the destination allocation protocol though maximum throughput is less even at $b = N$.

*Allocation Free Protocol*

The time division allocation cycle is collapsed into a single slot and any node with a message ready to send can transmit on a randomly chosen channel. This corresponds to the granting of $N \times N-1$ permissions in each slot. The system operation is thus completely random access. Minimum average queueing delays is shorter for this scheme but maximum throughput is also less.

In these protocol schemes there is no explicit mention as to how the destination receivers decide, or are directed, to tune to a particular frequency. In the fixed assignment source/destination allocation protocol and in the destination allocation protocol the frequency channel is specified for the sources selected. In the other schemes where multiple sources (more than the number of channels $b$) can be accessed from the permissions granted in a particular slot, it is not clear how receiver tuning is decided.

4.23) *Dynamic Time-Wavelength Division Multi-Access (DT-WDMA) protocol [CDR 90]*

As in other optical communication access control methods, this scheme is based on the use of multiple communication channels. In this scheme, however, contention does not occur for the use of communication channels but rather occurs at the receivers to which packets are destined. The network is again based on a star topology with stations employing fixed wavelength transmitters and tunable optical receivers. Each station therefore has its own assigned frequency that is uses for packet transmissions. A dedicated control channel is used for stations to broadcast their transmission intent so that the appropriate receivers can tune in to the frequency of the particular transmitters. Two transmitters are employed by each station; one for control channel transmissions, the other for the assigned frequency channel of the station. Access to the control channel is achieved on non-contention TDMA basis with fixed minislots

assigned to each network station. A global but distributed algorithm is used to resolve the contention of many transmitters wishing to access a given receiver. No explicit acknowledgements are thus required as each station is independently able to evaluate the result of ensuing contentions.

Stations with packets to transmit place a control packet in their assigned minislot on the control channel. This control packet contains the destination address, the delay experienced by the packet since its arrival at the station, and the mode of the transmitter. The mode information is included as both packet and circuit mode operation is possible with this protocol. A common clock is used in the network and allows packet delay to be calculated based on a common reference. Delay is the measure used to distributedly resolve the occurrence of packet conflicts at a receiver. The packet with the longest indicated delay is the one accepted at a receiver in the event of multiple control packets indicating a desire to communicate with a particular receiver. The delay time is measured as the sum of the delay experience until the transmission instant, and the time taken by the packet to reach the hub from the transmitting station. The delay time to the hub is dependent on the location of the station and is measured when the station enters the network.

Time is divided into slots of equal length. Both the slots on the data channels and the common channel are of the same length. The control channel slots are however further divided into the TDMA minislots assigned to individual stations. Stations having transmitted a control packet follow it with the transmission of the corresponding data packet in the next data slot. When several packets are addressed to the same receiver, the receiver can tune to one of the wavelengths and receiver the packet without interference from the other stations. Since all stations receive all min-slots, each transmitter can determine the outcome of its transmission without explicit acknowledgements. The central hub reduces the time between retransmissions since the outcome of packet transmissions are known in one roundtrip propagation time between the transmitter and the hub rather than between the transmitter and the receiver.

This protocol is able to offer very good throughput delay performance achieving a maximum throughput of approximately 0.6 even for very large numbers of users. There is no need for explicit acknowledgements since all stations monitor the common channel and the outcome of transmissions is known after a propagation delay equal to the roundtrip delay from a station to the hub. This keeps the average packet latency small irrespective of the number of stations in the network. The protocol operation is decentralized and hybrid switching is supported. The scheme is also stable as a receiver is able to correctly receive one packet even in the event of contention.

From the early optical protocol schemes we have examined, the elements of satellite multiaccess approaches are evident. The methods employed, however, have been essentially random access in operation and not involving the complexities of coordinated access based on distributed information. While this may simply represent the first level of protocol designs there may be certain inherent operational constraints that place limitations on the development of coordinated schemes. This we will seek to determine in the further course of our study.

The issues raised by the schemes presented above and those involved in the general design of optical network protocols are important in providing a better understanding of the design possibilities for multiaccess optical protocols. We shall thus discuss these issues in more detail before moving on our newly proposed scheme.

## 4.3 Issues in the design of optical network protocols

### 4.31 Network Model

The design of access protocols for optical networks presents unique challenges. The existence of multiple parallel communication channels requires a high degree of coordination in addition to the conventional multiaccess concerns of single channel communication. The access problem is now two-dimensional. Not only is there a requirement to regulate access on a given channel to ensure efficient usage of that communication resource, but there is also the need for an appropriate mechanism of channel assignment among users if effective overall bandwidth utilization is to be achieved. If the potential of optical fibre communication is to be fully harnessed, the efficient use of interchangeable frequency communication paths must be pursued. While fixed assignment of frequency channels can be applied, as in the scheme proposed by Chen *et al* [CDR 90], the ability to more appropriately utilize the communication resource in a bursty environment is lost without also having demand assignment of the frequency channels.

The single hop multiple channel environment being considered in the WDM approach to optical communication means that stations must be capable of communicating on any of the parallel frequency channels which are derived from a single mode fiber. As we pointed out earlier, this requires transceivers operating at each of the existing communication frequencies. Because of the prohibitive cost of this duplication of hardware, we seek network solutions which involve a minimum hardware requirement. This involves the use of tunable components.

### 4.32 Propagation Delay Concerns

The round trip propagation delay is one of the most significant factors affecting protocol design. It forces a distributed approach in order to avoid an additional hop of delay associated with centrally controlled schemes. With a network diameter of 100 km, the propagation delay time is about 500 µs. With a frequency channel capacity of 200 MHz and packets of length 1000 bytes, the packet

transmission time is 40 μs. Therefore the number of packet transmissions possible within a propagation delay interval is about 12. This is on the order of that encountered in satellite systems.

Unlike satellite networks in which propagation delay is generally fixed, the propagation delay of terrestrial systems is dependent on the diameter of the network. Therefore, for a fixed communication channel bandwidth, the number of packet transmissions per round trip delay increases with the size of the network. Protocol design must therefore be considered for fixed categories of network sizes. Alternately, to maintain the ratio of number of packet transmissions per round trip, the bandwidth of the communication channels must be reduced as the size of the network increases. This can be achieved by accordingly increasing the number of available channels (as detailed in the work by Chlamtac and Ganz [CG 88c] where an optimal bandwidth division can be found which maximizes system throughput and minimizes average packet delay).

## 4.33 Information Collection and Maintenance Demands

From our study of access schemes in the single channel satellite network environment, a number of effective schemes emerged which were able to achieve good throughput-delay characteristics by integrating random access and reservation access approaches. In such schemes, where distributed control was employed, one of the main operating requirements was the need for network users to maintain global queues which recorded the state of the channel for at least one frame duration (the frame duration being at least one round trip propagation delay period). In any attempt to introduce coordinated access schemes to the multichannel optical network, the information maintenance requirement is multiplied by the number of available frequency channels. In order to have reserved access to any frequency channel, or to combine random and reservation access, channel state information for each channel must be continuously available to all users.

Collecting channel state information to determine random access and reservation periods can be achieved in two ways. One approach is for each station to have receivers that continuously monitor

each frequency channel and maintain a record of the activity on the channel (for one frame duration) as in the single channel satellite case. This of course becomes an impractical solution for any significant number of optical frequency channels. The other alternative is to use the control channel as the medium for relaying all information necessary for the coordination of overall system access. In this case it is only necessary to continuously monitor the control channel. This however only addresses the problem of gathering the channel status information; limiting the amount of information that must be maintained for coordinated access is also of practical concern.

One approach to reducing the amount of information to be maintained for coordinated access schemes is the moving away from the slot-by-slot protocol operation to a multislot or circuit based operation. The number of data packet transmission periods per round trip delay will be determined by the length of allocated multislots. By having fewer transmission periods per round trip delay the data requirement for maintaining channel status will be reduced. This is discussed further in the next section as it also has the potential to effect a reduction of the traffic demand placed on the control channel.

## 4.34 Traffic Demands on the Control Channel

In addition to the data maintenance overhead in the use of coordinated access schemes, there is also the concern of the very large traffic demands imposed on the control channel by the need to request channel and slot assignments. As the data transmission requirements of network users continue to grow, the frequency of packet arrivals accordingly increases. This therefore results in a very high offered traffic load on the control channel. Simple calculations show that for a network of a few hundred broadband users, the packet arrival rate leads to a prohibitive processing requirement for control packets in a system operating on a slot-by-slot basis. The means to dealing with this situation lies in reducing the amount of information that must flow over the control channel without introducing flow control bottlenecks.

A ready means of reducing the potential overload on the control channel is use a contention-free TDMA approach. Stations would have fixed assigned slots on the control channel in which channel notifications can be made. Each time a station's TDMA slot is used, the reservation (notification) for all presently available packets can be made. Since the size of control packets is small relative to data packets and a sufficiently large bandwidth can be provided for the control channel, the average delay in accessing a TDMA reservation slot will be reasonable. The guarantee against collision may thus obtained at a small cost in average delay. The use of contention-free TDMA on the control channel also eliminates any instability problems that could result from excessive packet arrival rates, as would be associated with broadband communication.

Another means of addressing the problem of overload on the control channel is through the use of a circuit-based transmission approach. Under such a scheme, allocation of the frequency communication (sub)channels is provided on a virtual circuit basis rather than on a slot-by-slot basis. Under such a system, the arrival of a packet at a station initiates the submitting of a single notification (reservation) on the control channel. Then, once access to a data transmission channel has been set-up, all packets at the station are transmitted until the packet queue is empty. It is also be possible to place a limit on the number of packets that are transmitted once a data channel access has been initiated. This would provide for a multislot (maxislot) allocation system rather than a full circuit approach. The justification for such a variation lies in the ability to regulate access to the data transmission channels on a fixed interval basis. This would avoid the need to obtain explicit acknowledgements of the termination of use of a given channel.

This move to a circuit-based transmission allocation system also has the advantage of significantly reducing the data maintenance overhead associated with coordinated multichannel access. As mentioned previously, the granting of data channel usage on a slot-by-slot basis requires an excessive amount of information to be maintained by each network station. By coordinating and granting usage on a circuit (or multislot) basis, the problem of having to maintain large amounts of global queue

information for coordinated access of the multiple frequency channels can be reduced. With circuit assignment, it would be practically possible for individual stations to maintain channel usage information for a large number of data channels. Any use of circuit-based assignments would need to be limited since it is not compatible with bursty data transmissions.

A less desirable alternative is the setting-up of multiple control channels. Each control channel would have to be associated with a fixed group of data transmission frequency (sub)channels. Since users are interconnected over the control channel, this fixed partitioning of the system would be necessary to maintain overall coordination through the coordination of subsets of users. Given the need to have stations tuned to the control channel in order to identify intended transmissions, the assignment of users to a particular control channel would need to relatively stable. In an environment where users were being constantly assigned to new control channels there would be a need for a great amount of speed and traffic measurement overhead in the mechanism controlling the assignments (assumedly on a basis that reduces multiaccess contention amongst groups of users assigned to a single control channel).

In addition to coordination processing overhead, the time factor involved in relaying the assignment information and the time taken for channel frequencies to be tuned, would result in further inefficiencies. This means therefore that the assignment of groups of users to a particular control channel would have to be very slowly adaptive, or more probably, done on a fixed basis according to periodic measures of overall control channel traffic. With a fixed assignment there is nonetheless the inefficiency which may result from some groups experiencing high contention while others were relatively in active; the classic drawback of fixed assignment in a multiaccess environment. Therefore, this option for lightening the offered control channel traffic by introducing multiple control channels, creates a degree of sub-optimality in the bandwidth utilization of the system.

## 4.35 Limitations to Integrated Random Access/Reservation Approaches

Our earlier evaluation of multiaccess schemes showed the relative superiority of protocols that integrated random access and reservation mechanisms in their operation. This therefore encourages the consideration of a similar approach in the optical arena. The existence of multiple channels, however, presents a significant difficulty not previously encountered. In any scheme involving the coordinated use of the frequency bandwidth, stations must have some means of monitoring the activity on each of the communication channels. We partially addressed this problem by having all information notifying of transmissions to be sent on the control channel because of the practical infeasibility of having stations equipped with receivers to monitor each frequency channel simultaneously.

The principle of successful integrated schemes is to use contention access in certain allowable periods and resort to making reservations when not possible or in the event of collisions. Because of the need to notify intended receivers of the frequency channel to be used, random access data transmissions cannot precede reservations in an optical network. Any integrated scheme must first reserve (or notify of transmission) then attempt to reduce the reservation delay by means of random access transmission prior to the assigned reservation period. Hardware duplication constraints limit stations to only monitoring the control channel and the channel on which a data transmission is made. As a result the outcome of contention access is known only to the stations that attempted the transmissions on the particular channel within a given slot period. Since reservations had previously been made for these transmissions and may have been (or will be) received globally, a successful random access transmission will still have a reserved transmission period in a following frame unless some broadcast notification is made to relinquish the reservation.

The outcome of the contention transmissions themselves are however only known after a roundtrip delay and so any attempt to relinquish a reservation period by means of a notification on the control channel may only be received by other network stations after the reservation period has passed. As a result the transmission period may go unused and is not reassigned. In addition to being untimely, any

attempt to reassign reserved slots on the basis of successful contention transmissions will result in a reduced utilization of the control channel for the making of actual reservations. Thus without the ability of each station to continually monitor all frequency channels and directly maintain their own channel status information and reservation queues, the possibility of adequately integrating random and reservation access is severely curtailed.

## 4.36 Constraints on Reservation Approaches

Also of consideration in the design of optical protocol schemes is the employing of a reservation based approach. In the multichannel optical system each arriving packet a station will need to make reservations for transmission on one of the many frequency channels. It will not be possible, however, to group reservation requests for different destinations, since receivers must be individually notified of intended transmissions. With many packet arriving for different destination receivers, the number of separate reservations to be made will accordingly increase. Thus, in contrast to the single channel satellite case, it will not be possible to reduce the average waiting time to transmission by making reservations for groups of packets. A potential solution would be to increase the information content of control (reservation) packets to indicate multiple destination intent. Implementing such a scheme add further complexity and increase the operating overhead.

If the control channel is accessed on a contention basis, the time delay for reserved data transmissions may become greater than one propagation delay if the time between reservation packets for any station is less than a data packet transmission period. This can result in multiple reservation requests being made, with each subsequent transmission reserved for one data period later than the previous one, even though all reservations are made within one data time period. Thus as more successive reservation requests (control packets) are sent out by the given station, the time delay to the reserved transmission period for the data packets becomes greater than simply one propagation delay period. With contention access on the control channel the same problem can occur if a station is successful in transmitting multiple reservation requests within a time interval less than a data transmission period.

In the case of multislot reservations (with longer reserved periods), the problem becomes even more severe leading to much longer delays. The problems of longer delays may be further exacerbated if each station only possess a single transmitter since an additional time constraint is placed on when reservation requests can be transmitted.

The need for sequential assignment of transmission periods at a given station also leads to further complexity in the allocating of reservation slots on any frequency channel, among requesting stations. The reservation periods for different stations will now have to be interwoven with each other given the sequential transmission constraints at each station. Instead of simply making reservations for numbers of slots, transmission time periods will have to be reserved. This would necessitate global synchronization and the sending of transmission time information. In addition to this further difficulty, unless the reservation periods are of fixed length the resulting fragmentation of the transmission time on the data channels will result in reduced potential bandwidth utilization as well as longer delays.

We have seen in the previous section (4.2 above) a network model using fixed frequency transmitters and tunable receivers [CDR 90]. An alternate approach that may be considered is the use of tunable frequency transmitters with fixed receivers. In order to avoid collisions between multiple stations transmitting on a particular frequency (to communicate with a single receiver), there needs to be some means of access coordination. This implies some form of reservation scheme with stations notifying of their transmission intent. Because of the problems just described in such a reservation approach, the tunable transmitter/fixed frequency receiver model will also be limited in its performance efficiency if network stations are hardware limited to a single transmitter (or even two).

## 4.37 Conclusions on Issues of Optical Protocol Design

The framework of optical protocol design is determined by: the need to minimize the complexity and replication of transmit/receive mechanisms at each node, as dictated by cost and performance considerations; the necessity of significant data maintenance for coordinated operation; the desire to

manage the multiple alternate paths with minimum real time processing overhead; and the limitation posed by the channel scan time needed for determining the availability of idle channels [CG 88b], which dictates the maintenance of channel status information. These factors and the associated penalties must be addressed in the design of protocols aimed at achieving good system performance. Therefore, due to the greater complexity of the access coordination task, a greater level of complexity is involved than that encountered in satellite schemes.

It should also be pointed out that in optical networks given the large transmission bandwidths available, emphasis is better placed on reducing the average message delay and in achieving greater simplicity of operation rather than on bandwidth utilization efficiency through higher throughput capability. This of course may change as networks become larger, connecting greater numbers of users and as electronic communication bottlenecks are removed.

The issues raised above point to certain limitations and constraints that must be recognized in designing a practical access protocol for multichannel optical architectures. In the next chapter we present a protocol method that works within these constraints in seeking to maximize the utilization of vast communication resource of optical networks.

# Chapter 5

# Proposal for a New Optical Network Protocol

## 5.1 Introduction

In the previous chapter we examined many of the constraints and practical problems that affect the design of multiaccess protocols in optical communication environment. While wavelength division multiplexing (WDM) offers a way to utilize the vast communication bandwidths of single mode optical fibers, the difficulties inherent in the coordination of the multiple channel multiaccess communication places limitations on the type of protocol schemes that can be effectively employed. As seen in the schemes so far proposed these difficulties tend to encourage the development of random access based approaches.

In all such contention methods, retransmission exacts an additional penalty in bandwidth utilization as well as further delays as the number of collisions increases. In our satellite protocol study the purer random access schemes represented a simple, though limited, communication approach. That study also showed the superior throughput-delay performance of integrated random access/reservation schemes. The difficulty however of extending such integrated approaches to the multiple channel optical networks leads us to examine other means of achieving similar operating success. Desired is a means of reducing as much as possible the coordination overhead of any employed protocol scheme.

The scheme that we thus propose is a contention based approach which is derived from the URN protocol designed by Kleinrock and Yemini [KY 78] for multiaccess satellite communication (see section 3.401, p. 55). This adaptive protocol was shown to provide good throughput-delay performance over a wide range of load conditions and was relatively simple in its operation. It does not suffer the inherently large minimum average delays at low throughput associated with reservation schemes nor the complications operating in a multichannel environment. This approach also avoids the complexity and limitations of random access/reservation options as well as the high channel data maintenance requirements of coordinated optical schemes.

The protocol scheme operates in a grouped random access fashion by making selections of groups of users who are then granted permission to randomly transmit on a specified data channel. The principle of operation of the URN scheme is to determine the size of groups that must be formed in each transmission period such that in each group granted transmission permissions there is the probability that only one user has a packet to transmit. This relies on being able to make some determination of the number of users in the network who have packets to transmit at any given time. In extending the scheme to the multichannel optical network a separate control channel is also provided to allow receivers to be notified on intended transmissions. This extended URN protocol is shown to provide a good practical design for a multiaccess, WDM optical network.

## 5.2 Network architecture

The network model assumed in this newly proposed scheme is similar to that considered in the high-speed optical network protocols based on a WDM approach and employing a star topology. The network stations are interconnected via the broadcast medium of a passive star coupler (see figure 5.2.1 below). It is assumed that each network station is equipped with two optical transmitters and two receivers. Of these, one transmitter and one receiver is fixed and the other tunable. The fixed transmitter and

receiver are dedicated to the control channel. This not only represents a realistic and economic hardware provision but also leads to a truly multiaccess multichannel environment in which each station can communicate on any frequency channel with any other station in a single hop[15]. This network model also allows for the consideration of a communication environment in which the number of network nodes exceeds the number of available communication channels unlike the fixed frequency approach of Chen *et al* [CDR 90] (section 4.23; though it could be achieved if there is arbitration at the transmitters).

STAR
COUPLER



**Figure 5.2.1**     Star-configured multichannel optical network model.

The communication bandwidth of the optical fiber is divided up into $b + 1$ frequency channels each using a different wavelength. The $b$ channels operating at wavelengths $\lambda_1, \lambda_2, \ldots, \lambda_b$ are assigned to actual data traffic with one channel operating at wavelength $\lambda_0$ for control traffic [HKS 87]. There are $M$ stations in the network which can transmit or receive on any of the the data channels as well as communicate on the control channel.

With this multiple channel architecture each station must be informed of the channel it must listen to in order to receive intended transmissions and possibly the identity of the transmitting station. The control channel provides this function. Since network stations are limited to two receivers, it is not

---

15     Acampora [Aca 87], and others, have proposed network schemes in which stations have fixed transmission frequencies as well as fixed frequency receivers. Multihop communication between stations is therefore employed. The advantage however is that fast tunable components are not required.

possible for a station to monitor all the data channels and so all information relating to the use of the data channels must be sent over the control channel. One receiver is thus constantly tuned to the control channel while the other is sufficiently agile and ranging over the other $b$ channel wavelengths. Each station is equipped with a single buffer into which arriving packets are queued for transmission.

Each station also possesses two transmitters one of which is fixed and operates at the wavelength of the control channel, $\lambda_0$. The two transmitters allow for control information to be transmitted at the same time that a data transmission is taking place. This improves the operational efficiency of the protocol scheme over the use of a single transmitter which alternates between control and data channels.

.

The protocol operation is synchronous with time slotted into fixed length periods. This period is equal to the time taken for the transmission of a data packet on any of the parallel data channels. Time on the control channel is similarly slotted with operation on on the channels (including the control channel) being synchronized. Each time slot on the control channel is divided into $b$ minislots associated with each of the $b$ data channels. The control packets sent in these minislots consist of the transmitter address bits and the user address bits. The control packet is thus much smaller in length than the actual data packet, in our case $b$ times smaller. Because each of these minislots correspond to a particular data channel, the wavelength to be used in the communication between the transmitter and the receiver is already understood.

A station receiver monitoring the control channel is notified of an intended transmission by recognizing its address within a control packet. The second station receiver then tunes to the transmission wavelength for actual data packets. Since a receiver's address may appear in more than one control packet, indicating the intent by multiple stations to communicate with that receiver, the receiver simply tunes to the wavelength associated with the first control packet received which contains its address. This resolves destination conflicts. Since the control channel information for a given slot is

known one round trip period later, each station is able to determine at that time whether its packet was dropped because the intended receiver was tuned to another wavelength (one associated with a control packet appearing earlier in the reservation slot period).

## 5.3 The Extended URN (XURN) protocol scheme

This protocol is an extending of the adaptive URN scheme to meet the demands of the WDM optical communication environment. The essence of the new protocol's operation is based on the mechanisms devised by Kleinrock and Yemini [KY 78] to deal with the multiaccess problem on a single broadcast communication channel. The scheme is modified and augmented to operate in a multichannel network.

### 5.31 Details of the URN protocol

The goal of most adaptive control algorithms is to achieve a channel traffic rate, $G$, of 1. The main difficulty however is that since users are distributed, the number $n$ of ready users (that is, users with packets to transmit) at any time is generally not known to individual users.

Consider the case in which n is known to the individual users. (How n can be estimated will be discussed shortly.) An adaptive strategy for achieving the $G = 1$ condition is to have each ready user transmit in the next slot period with probability $1/n$. Kleinrock and Yemini [KY 78] proposed an alternate "pure" strategy for ready users to determine whether or not to transmit in the next slot: the probability of transmission is either 1 or 0. That is, some users have full channel access rights while other have none. A station that has been granted channel access and is also ready transmits in the next time slot.

The URN protocol is described using the following URN model. Each user is considered a colored ball in an URN; black for those ready and white for not ready (no packet to transmit). The access protocol is essentially a rule to sample balls from the URN. Let $k$ be the number of balls drawn from the URN. The

probability of a successful transmission (throughput) is that of getting exactly one black ball in the sample. This probability is given by an element of the Hypergeometric distribution:

$$\frac{\binom{k}{1}\binom{M-k}{n-1}}{\binom{M}{n}}$$

where $M$ is the total number of balls and $n$ the number of black balls. This expression is maximized when

$$k = [M/n] \qquad \text{where } [x] \text{ is the integer part of } x.$$

Not only does this value of $k$ maximize the probability of selecting exactly one black ball, but it also gives that the average number of black balls selected is equal to one (i.e $G = 1$).

. The URN protocol adapts smoothly to network load fluctuations. When the network is lightly loaded, a large number of users get channel access rights. For example when $n = 1$, then $k = M$; all users get full access rights, but only one (the single ready user) is going to make use of it. As the network load increases, $n$ increases and the number of users given access rights is reduced. When $n > M/2$ then $k = 1$ and the URN protocol becomes effectively a TDMA protocol (which is most suitable for a heavy load). The maximum channel throughput of URN is thus unity.

In the operation of this protocol scheme the two central issues of implementation are the obtaining of current values of $n$ by the individual stations, and the distributed coordination and assignment of the $k$ users given access rights in each slot.

Kleinrock and Yemini proposed the use of a binary erasure reservation subchannel for estimating $n$. An idle station which becomes ready ($n$ increases by 1) sends a message of a few bits in the subchannel. When a ready station becomes idle ($n$ decreases by 1), the condition is detected by other stations from examining its last packet or through positive acknowledgement. A collision on the reservation subchannel indicates that two or more users became ready. Since the probability of more than two users becoming ready at the same time was shown to be negligible, the error in determining $n$ in the event of

collision was negligible. The estimate of $n$ would be corrected whenever the network goes idle ($n = 0$). (This scheme was presented in the context of a zero propagation delay broadcast channel). Other heuristic algorithms for estimating $n$ in long propagation delay (satellite) channels were studied by Lam and Kleinrock [LK 75]. These are based on channel history (i.e., empty slots, successful transmissions, or collisions) being available to all channel users.

Implementation of the URN scheme also requires that individual stations agree upon $k$, the number of stations granted access rights, as well as their identity. This selection of the $k$ stations may done through the use of identical pseudorandom number generators at the individual stations, via a window mechanism, or other methods.

## 5.32 Operation of the XURN protocol

In this extended URN protocol, permissions are granted to groups of stations to transmit in a given slot and on a specified frequency channel. The information used to decide the granting of access rights similarly consists of the total number of busy stations in the network at the beginning of each given slot.

Because of the need to also inform receivers of intended transmissions in the multichannel environment, a separate control channel is maintained. The information sent on this control channel must also be conveyed prior to the actual data transmissions themselves. Thus, the control channel information is transmitted one slot period prior to the associated data packet transmissions. The control channel is slotted in time with each time slot period being equal is duration to the packet transmission time on the data channels. These time slots are further subdivided into $b$ minislots, where $b$ is the number of data frequency channels. Access to the control channel minislots is allocated in the same manner as the granting of transmission permissions to groups of stations. Members of each group of stations granted permission to transmit on a particular data channel are allowed to access a single minislot to convey the control information associated with their transmissions.

As in the case of the single channel URN scheme, an estimate of the number of ready stations must be available throughout the network. From this estimate, a determination is made of $k$, the number of stations to be granted access permissions on a single frequency channel, where $k = [M/n]$. However, because there are $b$ data channels in the network, the number of permissions granted is $bk$, with $k$ on each channel. As in the URN scheme, this approach maximizes the probability of providing access permission to exactly one ready station per data channel. If the number of ready stations is greater than the number of data channels, $n > b$, then the number of stations granted transmission permissions in a given slot is less than the total $M$. If, however, $n < b$ then all stations are granted transmission permission. The number of stations that would be assigned to each data frequency would be the total number, divided by the number of network channels. That is, $k = [M/b]$ permissions per channel.

The $k$ users assigned to each group can be randomly selected from among the total population of network stations. It is also possible to employ a windowing mechanism to determining which stations are granted access permissions in each slot time. This window will have $b$ sections each of size $k$ stations (see figure 5.32.1 below). All stations can be ordered according to their numbers along an imaginary circle. The $bk$ stations are selected by the window which rotates around the circle. After each transmission slot, the tail of the window is advanced along the circle to the head of the previous window position and the window size is once again set equal to $bk$ (based on the current estimate of $n$).

| $k = M/n$ | $k = M/n$ | $k = M/n$ | $\bullet\ \bullet\ \bullet$ | $k = M/n$ |
|-----------|-----------|-----------|-----------------------------|-----------|
| group 1   | group 2   | group 3   |                             | group b   |

**Figure 5.32.1**    Form for Windowing mechanism for XURN scheme.

As stated, each control channel time slot is divided into $b$ TDMA minislots. The stations that are assigned transmission permissions for access to the given data channels are the same ones allowed access to each of the associated minislots. For example, the stations that have been granted permission to access data channel 1 during the $(r+1)$ th time slot, are allowed to transmit control packets in

minislot 1 of the $r$ th time slot (see figure 5.32.2 below). Just as stations have been granted permissions such that there is maximum probability of only one station transmitting on a given channel within a time slot, so too the packet sent in the associated control minislot should correspond to that from the one ready station. Thus, successful transmission of a control packet on the control channel will correspond with a successful transmission on the associated data channel (with a probability of 1) and collision on the control channel will accordingly correspond with a collision on the data channels.

A receiver recognizing its address on the control channel immediately tunes to the wavelength specified by the control packet (using its second receiver). Receivers thus tune to the wavelength of the first control packet recognized with its address. With such an operation, multiple requests for a single receiver are readily resolved. As a result, however, it may be necessary to rotate the order of assignment of the TDMA minislots to specific data channels. By so doing there is no asymmetry in the resolving of destination conflicts that could favor stations assigned to lower numbered minislots. (This is not however a significant concern given the variability of the window size as it advances around the circle of users). The figure below shows the timing diagram of the operation of XURN scheme.



Figure 5.32.2    Timing diagram of the operation of XURN scheme.

103

By continuously monitoring the activity on the control channel (using their fixed receiver) each station will have the information necessary to implement an algorithm for estimating $n$. The network (all channels) history, comprising empty slots, successful transmissions, and collisions, will be available from monitoring the control channel. While single channel heuristic algorithms for estimating $n$ have been studied [LK 75, Bertsekas and Gallager 87], it will be necessary to study the design of similar algorithms for the multichannel case. However, since all the necessary information is available, it is felt that similar algorithms can be devised.

A number of the features of this XURN scheme makes it particularly suited to operation in the multichannel optical environment and allows it to overcome some of the operational constraints described previously in section 4.3. The random access basis of the protocol eliminates the need for the high data maintenance necessary in coordinated approaches (though that information may still be necessary for the estimation of $n$).

This protocol also employs only a single queue at each network station and therefore only requires a single buffer. The use of a single buffer at each station, for all incoming packets, reduces the hardware requirement at each station. This is instead of the necessity for a separate buffer for each destination station as is necessary for some schemes.

Another problem reduced is the potential for traffic overload on the control channel. The number of control packets handled per data slot period is fixed at $b$. For example, a 1000 byte packet has a transmission time of 40 μs on a 200 b/s data channel. In such a case, the number of control packets to be processed would be fixed at $2.5 \times 10^6$ per second for a network of 100 channels.

Additionally, the scheme does not suffer the potentially mounting delays in packet transmission associated with reservation schemes, which (as explained in section 4.36) may occur as a result of the limited number of transmitters per station. The XURN scheme is instead able to capture the relatively

simple operation as well as adaptive property of its URN basis. Also, since groups of stations are assigned to a given frequency channel, the scheme is not highly susceptible to errors in the estimate of $n$. This is important for a distributedly controlled protocol.

## 5.4 Analysis of the XURN protocol scheme

In our analysis of the XURN scheme we shall consider the case in which each network station has only a single-packet buffer. This case is more directly solved and provides for a clearer understanding of the operation of the XURN protocol. We shall later consider the case of stations having multiple packet buffers and briefly mention the performance differences that result, and the most effective methods devised for conducting the analysis. For analytical convenience we also consider the case of stations being randomly assigned transmission permissions instead of using a windowing scheme.

### 5.41 XURN scheme for network stations with single-packet buffers.

The operation of the XURN scheme can be described by a discrete-time Markov chain. Let $n$ be the number of ready stations at the beginning of a given slot. The Markov chain is finite with the number of states being equal to $M$, the number of network stations. The time unit for the system is one slot. i.e, the transmission time of a data packet. The transition structure of this chain is shown in figure 5.41.1 below.



Figure 5.41.1    Markov chain for XURN protocol. The state (i.e, number of ready stations) can decrease by at most $b$ per transition, but can increase by an arbitrary amount.

Consistent with the positioning of the transmission window (see figure 5.31 above), each ready station which is granted permission to transmit in a given slot will do so. From the URN scheme, for $k$ stations granted transmission permission in a given slot, the probability of a successful transmission is given by:

$$\text{Probability of a successful transmission} = \frac{\binom{k}{1}\binom{M-k}{n-1}}{\binom{M}{n}} \tag{5.41.01}$$

where $k = \left[\frac{M}{n}\right]$ such as to maximize the probability of a successful transmission.

For the XURN scheme we now have $b$ groups of $k$ stations to which transmission permissions are granted since there are $b$ available data transmission channels. Let us now consider the probability of successful transmissions on these channels.

The probability of choosing exactly $b$ ready stations among $bk$ stations selected (to give transmission permissions) is given by the hypergeometric distribution:

$$\frac{\binom{bk}{b}\binom{M-bk}{n-b}}{\binom{M}{n}} \tag{5.41.02}$$

The random groups of stations are selected such that $k = \left[\frac{M}{n}\right]$ for $n > b$ and $k = \left[\frac{M}{b}\right]$ for $n \leq b$ stations.

Note: It is important to set these limits on the value for $k$, since if we maintain $k = \left[\frac{M}{n}\right]$ for $n \leq b$ then the number of stations granted transmission permissions in a given slot would be greater than $M$, the number of existing in the network.

We can get a sense of the performance of the protocol scheme by considering the expected number of successes achievable in a given slot period.

Let,

$\quad\quad P(r_i) = \text{Prob. } [r_i \text{ ready stations in the } i \text{ th group}]$

then,

$$P(r_1) = \frac{\binom{k}{r_1}\binom{bk-k}{n-r_1}}{\binom{bk}{n}} \quad \text{and}$$

$$P(r_2) = \frac{\binom{k}{r_2}\binom{bk-2k}{n-r_1-r_2}}{\binom{bk-k}{n-r_1}} \quad \text{... etc.}$$

Prob. [success in $i$ th group] = Prob. [$r_i = 1$]

therefore,

Expected number of successes in network in a given slot period

$$= \sum_{i=1}^{b} \text{Prob. } [r_i = 1]$$

Irrespective of conditionality,

Expectation of Sum = Sum of Expectation

Hence,

Expected number of successes in network in a given slot period

= Number of groups x Expectation of success in any one group

$$= b \times \frac{\binom{k}{1}\binom{M-k}{n-1}}{\binom{M}{n}} \tag{5.41.03}$$

where $M$, $n$, and $k$, are defined as above in equation (5.41.01).

This equation confirms that the choice of $k$ to maximize the expected number of successes in the system, is the same as that given for the URN scheme (with the additional limit specification given in equation 5.41.02). It is clearly seen therefore that the expected number of successful channel transmissions for the XURN scheme is multiplied by the number of available frequency channels.

The probability of a successful transmission among any one of the $b$ groups of $k$ stations is equal to the probability that there is exactly one ready station among the $k$ chosen for the group, given that the total number of stations selected (i.e given transmission permissions) is $bk$, and given that there are a total of $n'$ ready stations among these.

Prob. of a successful transmission within a single group of stations

$$= \sum_{n'=1}^{\min(bk,\, n)} \text{Prob [exactly one ready station within group} \mid n']$$

$$\times \text{ Prob } [n' \text{ ready stations among } bk \text{ selected}]$$

And,

$$\text{Prob. } [n' \text{ ready stations among } bk \text{ selected}] = \frac{\binom{bk}{n'}\binom{M-bk}{n-n'}}{\binom{M}{n}} \qquad (5.41.04)$$

where $k = \left[\dfrac{M}{n}\right]$ for $n > b$ and $k = \left[\dfrac{M}{b}\right]$ for $n \le b$.

Therefore,

$$\cdot \quad \text{Prob. [successful trans. in a group of } k \text{ stations]} = \sum_{n'=1}^{\min(bk,\, n)} \left[\frac{\binom{k}{1}\binom{bk-k}{n'-1}}{\binom{bk}{n'}}\right] \cdot \left[\frac{\binom{bk}{n'}\binom{M-bk}{n-n'}}{\binom{M}{n}}\right] \qquad (5.41.05)$$

Given that there are $n'$ ready stations among the $bk$ stations granted transmission permissions, the probability of a successful transmission in any one group is conditional on the probability of successes or failures in the other remaining $(b-1)$ groups.

That is,

Prob. of exactly $q$ successful channel transmissions among $b$ groups

$$= \sum_{n'=q}^{\min(bk,\, n)} \text{Prob. } [n' \text{ ready stations among } bk \text{ selected}]$$

$$\times \text{[Number of ways of arranging } q \text{ successes among } b \text{ groups]}$$

$$\times \text{ Prob. [successful trans. in one group} \mid n']$$

$$\times \text{ Prob. [successful trans. in a second group} \mid \text{successful trans. in one group, } n']$$

$$\times \text{ Prob. [successful trans in a third group} \mid \text{successful trans. in two groups, } n']$$

$$\cdots \cdots \cdots \cdots$$

$$\times \text{ Prob. [successful trans. in a } q\text{th group} \mid \text{successful trans. in } (q-1) \text{ groups, } n']$$

$$\times \text{ Prob. [no successes in remaining } (b-q) \text{ groups of stations} \mid q \text{ successes, } n']$$

$$(5.41.06)$$

The limits of $n'$ extends from $q$ to $\min(bk, n)$ since the number of ready stations among the $bk$ granted transmission permissions is always greater than, or equal to, the possible number of successful transmissions. Also, $n'$ is always less than or equal to $\min(bk, n)$.

The dependencies in the above equation presents particular difficulties in deriving an exact analytical expression for the probability of multiple successful transmissions in the network, in a given time slot. We thus now examine the elements of this equation in more detail.

The order number of the groups of the in which successful transmissions occur does not matter and so we can consider the probability of successes in the first q groups selected. Using the hypergeometric given by equation 5.41.02 above, we can calculate the probability of successful transmissions in any particular group given the occurrence of successes in preceding groups. That is,

Prob. [successful trans. in first $q$ groups | $n'$]

$$
= \left[\frac{\binom{k}{1}\binom{bk-k}{n'-1}}{\binom{bk}{n'}}\right] \cdot \left[\frac{\binom{k}{1}\binom{bk-2k}{n'-2}}{\binom{bk-k}{n'-1}}\right] \cdot \left[\frac{\binom{k}{1}\binom{bk-3k}{n'-3}}{\binom{bk-2k}{n'-2}}\right] \cdot \left[\frac{\binom{k}{1}\binom{bk-4k}{n'-4}}{\binom{bk-3k}{n'-3}}\right] \cdots \cdots \left[\frac{\binom{k}{1}\binom{bk-qk}{n'-q}}{\binom{bk-(q-1)k}{n'-(q-1)}}\right]
$$

$$
= \left[\frac{\binom{k}{1}^q \binom{bk-qk}{n'-q}}{\binom{bk}{n'}}\right] \tag{5.41.07}
$$

The difficulty in arriving at an exact analytical expression for the original equation given on the page above stems from the need to calculate the probability of an exact number of failures, where the probability of failure in any group is also conditional on the occurrence of failures (and successes) in preceding groups. A failed transmission, however, results from the occurrence of zero, or more than one ready station in a given group of stations. Calculating the probability of failure in a given group conditional on events in previous groups can be approached in two basic ways.

Let us consider the probability of failure in consecutive groups where the total number of stations is $N$ (that is, $N/k$ groups) and the number of ready stations is known. We assume it is $n''$.

Prob. [failure in one groups | $n''$]

$$= \sum_{\substack{i=0 \\ i \neq 1}}^{\min(n'', k)} \frac{\binom{k}{i}\binom{N-k}{n''-i}}{\binom{N}{n''}}$$

Prob. [failure in second group | failure in first group, $n''$]

= Prob. [0 ready station in first group]

x { Prob. [0 ready stations in second group | 0 ready in first, $n''$]

+ Prob. [2 ready stations in second group | 0 ready in first group]

+ Prob. [3 ready stations in second group | 0 ready in first group] + ... }

+ Prob. [2 ready in first group]

x { Prob. [0 ready stations in second group | 2 ready in first, $n''$]

+ Prob. [2 ready stations in second group | 2 ready in first group]

+ Prob. [3 ready stations in second group | 2 ready in first group] + ... }

+ Prob. [3 ready in first group]

x { Prob. [0 ready stations in second group | 3 ready in first, $n''$]

+ Prob. [2 ready stations in second group | 3 ready in first group]

+ Prob. [3 ready stations in second group | 3 ready in first group] + ... }

etc.

$$= \sum_{\substack{j=0 \\ j \neq 1}}^{\min(n'', k)} \binom{k}{j} \sum_{\substack{i=0 \\ i \neq 1}}^{\min(n''-j, k)} \frac{\binom{k}{i}\binom{N-2k}{n''-j-i}}{\binom{N}{n''}}$$

This equation fast becomes intractable as the number of consecutive group failures increases. For each succeeding group we have to additionally consider the many permutations possible for in having 0, 2, 3, 4, ..., min($n''$, $k$), ready stations in the preceding ($q$-1) failed groups.

The other alternative to calculating the probability of multiple group failures is based on a count of the number of ways of arranging for all failures among groups of stations given that there is a known number of ready stations. That is,

Prob. [exactly $q$ failures in $q$ groups of stations | $n''$ total ready stations]

= [Number of ways of arranging for all failures with $n''$ ready stations among q groups]

divided by     [Total number of ways of arranging $qk$ ($q$ groups) stations]

Given that a failure is defined as the occurrence of zero or more than one ready station per groups we can examine the number of possible of ways failed groups can be arranged. Under unrestricted conditions it may be possible to derive analytic expressions for number of failure combinations and total number of station arrangements. In our case however two restrictions come into play. The first is the limit on the number of stations per group arrangement which must be $k$ for each group. This also affects the calculation of the denominator in the above expression. Additionally, the number of groups in the is fixed, and so the arrangement of failures are limited to combinations that extend across a maximum of q groups (in the above example). With these constraints deriving a general analytic expression becomes an intractable problem.

Because of the difficulties encountered in deriving exact expressions for analysis of the operation of the protocol the approaches that must be pursued involve either simulations or the application of computer programs that calculate the necessary probabilities by generating and counting failure and success arrangement possibilities. In the result section (5.5) later in this chapter, we use a similar counting approach to generate exact results for a small network case.

Using a computer program to count possible arrangements may not always be a satisfactory approach since the possible number of combinations can become excessive as the network size increases since it has a factorial dependence. It is also possible to consider the approximate operation of the scheme by

assuming independence between the groups of stations. This may be valid under conditions in which the number of busy stations in the network is very large. By removing the conditionality between groups a closed form analytic solution can be defined for the system operation. This issue is not further pursued at this time.

For the purpose of further analysis of the XURN protocol, let

$f(n, q)$ = Prob. of exactly $q$ successful channel transmissions given $n$ busy stations in network (as defined by equation 5.41.06)

This multichannel network architecture with single receiver stations, means that there is the possibility of multiple successful channel transmissions intended for a given receiver. This receiver is however only capable of correctly receiving one packet transmission per slot. Thus we consider a packet to be **completely** successful iff the following conditions are obeyed sequentially:

1)    the packet is successfully transmitted on a data channel with no channel collision, and

2)    the packet is successfully received (i.e, the associated control packet is the first recognized by the intended receiver).

Therefore, consider now the case of packet receptions. For each packet successfully transmitted on a data channel, there is an equal likelihood that it is destined for one of $(M-1)$ stations (where we assume a symmetric network of users). If we have $q$ successful channel transmissions in a slot period then these packets arriving at destination receivers correspond with successively choosing $q$ out of $(M-1)$ possible destinations (with repetitions being possible).

The operation of the XURN scheme is such that the first of multiple packets arriving at any receiver is always correctly received (i.e, with probability 1). Since the receiver tunes to the wavelength of the first packet recognized from information on the control channel, all other packets are dropped.

However, calculating the probability of packet receptions at a number of distinct stations is complicated by the fact that the $(M-1)$ destination choices for each station transmitter, is different (though there is a majority overlap for any two stations). Consider the following analysis for deriving the probability of a given number of correctly received packets.[16]

Let $g(q, s)$ be the probability that exactly $s$ packets will be successfully received in a slot given $q$ packets were successfully transmitted without channel collisions. To calculate $g(q, s)$ we define $\beta_u(M, q)$ as the probability that there is no packet to $u$ specific destinations among the $M$ stations, given that there were $q$ successful channel transmissions in the slot. Assuming that each destination is equally likely to be selected, with probability $\dfrac{1}{M-1}$ we obtain:

$$\beta_u(M, q) = \sum_{m=\max(0,\ q-M+u)}^{\min(q,\ u)} \text{Prob. (the source of } q\text{-}m \text{ packets is among } M\text{-}u \text{ stations and the source of } m$$

$$\text{packets is among } u \text{ stations)}$$

$$\text{x} \quad \text{Prob. (the destination of } q\text{-}m \text{ packets is among } M\text{-}u \text{ stations)}$$

$$\text{x} \quad \text{Prob. (the destination of } m \text{ packets is among } M\text{-}u \text{ stations)}$$

$$= \sum_{m=\max(0,\ q-M+u)}^{\min(q,\ u)} \frac{\binom{M-u}{q-m}\binom{u}{m}}{\binom{M}{q}} \cdot \left(\frac{M-u-1}{M-1}\right)^{q-m} \cdot \left(\frac{M-u}{M-1}\right)^{m} \tag{5.41.08}$$

Note that $g(q, s)$ is equivalent to the probability that there are exactly $M$-$s$ stations to which there is no destined packet. Consequently $g(q, s)$ can be computed using the Inclusion-Exclusion principle and $\beta_u(M, q)$ , to obtain:

$$g(q, s) = \begin{cases} \displaystyle\sum_{u=M-s}^{M} (-1)^{u-M+s}\binom{u}{M-s}\binom{M}{u}\beta u(M, q) & \text{for } s < b \\[2em] \displaystyle\sum_{v=b}^{M}\sum_{u=M-v}^{M} (-1)^{u-M+v}\binom{u}{M-v}\binom{M}{u}\beta u(M, q) & \text{for } s = b \end{cases} \tag{5.41.09}$$

---

[16]    This analysis is adapted from that given in paper by Ganz and Chlamtac [GC 87].

As stated above, the probability of a 'completely' successful packet transmission is dependent on both the probability of a successful channel transmission as well as the probability of a successful reception. Hence, within a given transmission slot, the probability of $s$ completely successful transmissions in the entire network (of $b$ channels) is given by:

$$\sum_{q=s}^{min(b,\,n)} \text{Prob. } [q \text{ successful channel transmissions}]$$

$$\times \text{ Prob } [s \text{ successful packet receptions } | \ q \text{ successful channel transmissions}]$$

that is,

Prob. of $s$ **completely** successful transmissions, given $n$ ready stations in network, is defined by:

$$P_t(n, s) = \sum_{q=s}^{min(b,\,n)} f(n, q) \cdot g(q, s) \tag{5.41.10}$$

where $f(n, q)$ is defined by expression 5.4105, and $g(q, s)$ is defined by equation 5.4109.

Let us now consider packet arrivals at the network stations. We assume that packet arrivals per transmission slot, into the network, is Poisson distributed with mean $\lambda$. Thus the mean arrival rate at each network station is $\lambda/M$. The probability of no packet arrivals at a network station within a slot period is $e^{-\lambda/M}$, hence:

Probability of packet arrivals at any station within a slot period $= 1 - e^{-\lambda/M}$. (5.41.11)

Note: As above, we assume a symmetric network of users.

From one slot to the next, the state (i.e, the number of ready stations) of our Markov chain increases according to the difference between the number of non-ready stations (i.e, stations with no packets awaiting transmission) to which new packets have arrived, and the number of ready stations from which there have been successfully transmitted packets. The number of non-ready stations in the

114

network is equal to $(M-n)$. Let $P_a(n, j)$ be the probability of $j$ packet arrivals at non-ready stations given there are $n$ ready stations in network.

It thus follows,

Prob. of $j$ new packet arrivals at non-ready stations, given $n$ ready stations in network is defined by:

$$P_a(n, j) = \binom{M-n}{j}\left(1 - e^{-\lambda/M}\right)^j \left(e^{-\lambda/M}\right)^{M-n-j} \tag{5.41.12}$$

Hence the transition probability of going from state $n$ to $n+i$ is given by:

$$P_{n,\,n+i} = \text{Prob. } [(j-s) = i] = \sum_{s=0}^{\min\,(b,n)} P_a\left(n, j = (i + s)\right) . P_t(n, s) \qquad \text{for } j \geq 0 \tag{5.41.13}$$

also note,

$$P_{n,\,n+i} = 0 \qquad \text{for } (n+i) - n < -b$$

where $j$ = number of non-ready stations at which there are packet arrivals, and

$s$ = number of completely successful transmissions (which includes successful receptions) among ready stations granted permissions.

From equations 5.41.10 and 5.41.12 we can calculate the probabilities $P_a(n, j)$ and $P_t(n, s)$ and determine the transition probabilities of the Markov chain that describes the operation of the XURN protocol.

The transition probability matrix of the Markov chain for the system can be defined as [P], such that the steady state equation for the distribution of the number of ready stations (i.e, stations with packets to transmit) has the form:

$$\pi = \pi \, [P] \tag{5.41.14}$$

115

where $\pi$ is the probability row vector $\pi = (\pi_1, \pi_2, \ldots, \pi_M)$ whose components are the probabilities of being in the various states in the limit as time goes to infinity. That is, $\pi_n$ is the steady state probability of finding $n$ ready stations in the system. Also,

$$\sum_{i=1}^{M} \pi_i = 1 \qquad (5.41.15)$$

Because the state of the Markov chain can decrease by at most $b$ per transition, the matrix $[P]$ is of form:

$$\begin{bmatrix} p_{00}, p_{01}, & \ldots & p_{0M} \\ p_{10}, p_{11}, & \ldots & p_{1M} \\ p_{20}, p_{21}, & \ldots & p_{2M} \\ \ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots \\ p_{b0}, p_{b1}, & \ldots & p_{bM} \\ \ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots \\ 0 & \ldots & \ldots \\ 0 & 0 \quad p_{Mb} \ldots & p_{MM} \end{bmatrix}$$

. For the operation of the protocol scheme as defined, the Markov chain has only one recurrent class and it is also aperiodic (i.e, it can be defined to be ergodic). As a result, theory indicates that each row of $[P]^n$ converges to the unique probability vector solution of equation 5.41.14 above. This probability vector solution is unique due to the existence of a single recurrent class in the chain. From a solution of the steady state equations we can compute the performance measures of throughput and average delay.

The system throughput is given by:

$$S = \sum_{n=1}^{M} \pi_n \sum_{s=1}^{\min(b,n)} s \cdot P_t(n, s) \qquad (5.41.16)$$

where $\pi_n$ is the steady state probability as defined above, $b$ is the maximum number of successful transmissions per slot (equal to the number of data channels), and $P_t(n,s)$ is the prob. that there are $s$ completely successful packet transmissions in a slot period given there are $n$ ready stations in the network at the beginning of the slot (as defined by equation 5.4110 above).

Using Little's theorem, the average packet delay is given by:

$$D = \frac{Q}{S} \qquad (5.41.17)$$

where Q is the average number of ready stations in the network and is given by:

$$Q = \sum_{n=1}^{M} n.\pi_n \qquad (5.41.18)$$

These equations provide the basis from which we can derive an assessment of the performance of the XURN scheme for the restricted case of stations having single packet buffers.

## 5.42 XURN scheme for network stations with multiple-packet buffers.

The case of stations with single packet buffers offers a more straightforward analytical exercise and leads to a better initial understanding of the operation of the proposed XURN protocol scheme. It does not however represent a sufficiently practical implementation since it assumes that packets arriving at a ready station must be dropped. We thus briefly consider the case in which each network station is equipped with a finite buffer capable of queueing $L$ *packets*. The relevant performance differences are that multiple-packet buffers allow for greater throughput values and lower rejection probabilities, though at the expense of an increase in delay.

It is possible to model such a system as a discrete-time queueing system consisting of M queues, each with a buffer capacity of $L$ packets, that interact according to the rules of the XURN protocol. The behavior can be described as an M-dimensional Markov chain with state vector $w = (w_1, w_2, \ldots, w_M)$ where $w_i$ is the number of packets in queue at the $i$th station, $i = 1, 2, \ldots, M$. Given the symmetry of the network stations, it is possible to describe the state of the system by the vector $n = (n_0, n_1, n_2, \ldots, n_L)$, where $n_i$ is the number of stations having $i$ packets in queue, $i = 0, 1, \ldots, L$. In this case $n_0 = M - n_1 - n_2 - \ldots - n_L$ and the number of possible states is reduced to $\binom{M+L-1}{L}$. Apart from certain relatively simple cases, there are no precise theoretical results for such buffered schemes [SKP 86]. Instead methods have been proposed for approximate analysis.

Sykas, Karvelas and Protonotarios [SKP 86] proposed a method of approximate analysis for such buffered multiple access schemes which significantly reduces the state space of the model and provides for a computationally more efficient means of analysis. Their approach involves using the number of busy (ready) stations in the system together with the queue length at a particular station to derive an efficient, though approximate, description of the protocol scheme being analyzed. A two-dimensional

Markov chain is constructed, now with a lower number of states $(L+1) \times M$, and (utilizing the symmetry of the system) the underlying transition probabilities are computed. In their paper they considered the application of this approach to the case of the basic URN scheme.

Ganz and Chlamtac [GC 87] further introduced an approximate analysis queueing models for finite buffer systems which also reduces the number of states in the systems' Markov chains to $(L+1) \times M$. However, as opposed to previous approximate solutions which involved the solution of non-linear equations, their approximation involves the solution of (a smaller number of) *linear* equations only. In addition, and more important from the point of view of the XURN scheme, the introduced analytical approach allows a model generalization that can be used to evaluate systems in which several successful transmissions can take place in parallel. We will not pursue the analysis of the multiple packet buffer case any further at this time. In the event it is desired, it can be efficiently performed as indicated in the above sources.

## 5.5 Results of XURN analysis

In this section we configure a small sized network in order to generate the delay throughput characteristic under the operation of the XURN protocol. We consider a network consisting of $M = 10$ stations connected across a star coupled WDM optical network with the number of available frequency channels, $b = 4$. Appendix A details the calculations performed to achieve the following delay-throughput results.



**Figure 5.5.1**  Average delay vs system throughput for XURN protocol with network of M = 10 stations and b = 4 wavelengths.

These results are very good and in keeping with that expected given the URN basis of the protocol. It thus confirms the ability to extend the success of that adaptive protocol to a multichannel optical environment. Comparing this performance with the results of the single channel case given in section 3.401 (p. 68), the XURN scheme as analyzed does not approach the maximum throughput of $b = 4$. This however is due to the model used for the analysis in which a packet, or packets, arriving at an already

busy station are simply dropped due to the single packet buffer assumed at each station. This prevents the expected TDMA type result from being fully realized at high system load.

This system throughput corresponds to a maximum normalized throughput of approximately 0.7. This is indeed a very good performance result given that the maximum throughput obtained for other optical schemes have been on the order of that achievable for ALOHA and slotted ALOHA type operation (the schemes by Habbab *et al* [HKS 87] that have achieved comparable throughput ranges have been based on the use of channel sensing, an approach unsuitable for large high speed networks).



Figure 5.5.2    System throughput vs Input rate (at each station) for XURN protocol.

In this plot of throughput versus input rate we seen the effects of the analytical model with the linearity of the curve rolling off some measure before the maximum potential throughput is reached. The characteristic is nonetheless as expected from the URN results originally obtained [KY 78]. These initial results obtained auger very well for the prospect of applying the XURN scheme to larger high speed optical networks.

## 5.6 Conclusions

In addition to the very good delay-throughput performance demonstrated, the XURN scheme offers a number of other operating advantages. Inherited from the URN scheme is the ability of the protocol to adapt to changes in network load conditions; changing from unrestricted ALOHA type access during low load conditions, to TDMA operation when the network is heavily loaded. As described, the scheme's hardware requirement is limited, needing only two transmitters and two receivers and a single buffer at each network station. The system is also totally flexible allowing each station the possibility to transmit on any frequency channel subject to the scheme's operating rules. The network is also truly multiaccess being able to accommodate a number of stations greater than the number of existing frequency channels. Another important practical feature of the protocol is that it is not highly susceptible to errors in the estimate of the number of busy stations, $n$. This could have led to potential problems in the distributed environment.

The reservation (notification) mechanism employed on the control channel resolves the potential overload problem that may occur in some of the proposed optical schemes if very large data rates are introduced in a network of many stations. The random access basis of the XURN scheme also avoids many of the coordination and synchronization problems that could affect a multichannel system. These many features thus add to producing a protocol whose overall performance characteristics are quite attractive for meeting the needs of high speed optical communication.

### 5.61 Further Work

A critical element in the operation of the XURN scheme is the determination of the number of busy stations, $n$, in the network at any given time. Throughout our discussions we have assumed that knowledge of this variable (or an estimate) is available at each station during every slot interval. While algorithms have been studied for calculating estimates of $n$ in the single channel case, it is necessary that similar methods be devised for the multichannel optical environment. Lam and

Kleinrock [LK 75] (and others [Gallager and Bertsekas, 87]) have studied the situation for satellite channels and have proposed methods based on channel history information. Since each network station in the XURN protocol constantly monitors the control channel, the channel history for all the network channels is obtainable. We thus conclude that it would be possible to distributedly calculate estimates of $n$ for the XURN scheme and recommend for further study the devising of a suitable algorithm.

Another area to be considered is the simulation approach required to generate channel success probabilities as described in section 5.41. Because of the intractability of defining exact analytical expressions it is necessary to construct simulation models that can allow the desire probabilities to be obtained for different system configurations.

# Chapter 6

# Understanding the Process of Technological Innovation

## 6.1 Introduction

As detailed in chapter 1, the approach adopted in this study of the design of a multiaccess communication protocol was undertaken on the basis of rational technical criteria. The problem of increasing the communication potential of optical networks from a basis of satellite communication was addressed from the point of the technical similarities that exists between established and understood satellite communication environment and the new emerging optical networks. By developing a further understanding of an older technology we have provided a basis for moving forward. Technically, the approach represented a logical path to achieving the design objectives. However, the effect of this starting logic may have constrained the results obtained. Understanding and highlighting issues of this nature and others surrounding the process of technological innovation and development will be the objective of this final stage of our study.

In this segment of our thesis work we thus take a step backward and objectively examine both the process and the product of our technical study. We seek to better understand the process of technological innovation and factors that affect the development of technologies. Of particular interest is the process involved in directed theoretical research undertakings. Unlike a number of the existing studies which

have looked at the process of technological innovation and adoption through an exploration of macro-level influences, we will focus our concerns on the micro-level elements that most affect the development of research-based studies. This focus thus relates more to the internal mechanisms that shape the course of technology.

We begin by providing a general definition of technological innovation which will establish our base reference for the discussions that follow. The process of innovation encompasses the use of knowledge for the generation and practical application of new and viable ideas [Holt, 87]. Innovation is thus a new or unique application of particular knowledge. The technology itself is also more than just the products, tools, or their usage, but includes the embodiment of certain knowledge and understanding. In this regard innovation is distinguished from invention as it goes beyond discovery, but rather represents an application. The definition that we apply is more broad than the purely economic or commercial view, employed in much of the innovation literature, that defines innovation as the bringing to market of a new product or process.

The overall process of technology innovation is often defined as comprising four main stages: discovery, invention, application and diffusion. This process is complex and often evolving over a long period of time. The various stages seldom occur in the same place, and usually different organizations and persons contribute to the process. Though not implying a rigid, linear process model, we nonetheless define technological innovation in this work in a more limited manner that focuses on the earlier stages of the technology development; namely the discovery and invention stages.

Technological innovation and development are brought about by a mixture of factors. The emphasis on the technical and economic determinants tends to obscure the influence of other more subtle contributors on the direction of emerging technologies. In the course of this chapter we will explore these other contributory factors based on the experiences derived from the study of protocols for high speed optical networks.

We shall begin by presenting the general framework in which technology development exists. This includes the organizational, socioeconomic, and political context that surrounds the development of technology. This is followed by an analysis of the specific factors that influence the internal process. From this enhanced understanding we recommend initiatives that will seek to ensure that the objectives of future technological studies can be more efficiently achieved.

## 6.2 The Context of Technological Innovation

In the field of "innovation process research" many different approaches have been taken in defining the context of technological innovation. Some researchers have focussed on macro variables such as tax, social and other governmental policies, or on micro variables such as characteristics of innovation adopters, and the role of individual researchers. Others have look at the organizational context in which the effects of these variables are played out. In this analysis, however, we take a more limited view and concentrate on the micro variables which make up the internal mechanisms of the technological innovation process; factors related to the design and problem solving stages of technological development. These variables however are not independent of the external macro level factors that shape development. We briefly examine this broader context so as to have a more complete view of the framework within which the process lies. The figure below provides a general overview of this larger context.

**Figure 6.1.1**   Overview of the technology development process

The main focus of our examination will be on the areas of technology introduction and the internal workings of the technology development process itself. In this regard we take somewhat the "internalist" view of the dynamic tendencies of the technology. We thus assume an evolution based on technological choices much akin to evolution in the Darwinian sense based on natural selection. The technical possibilities open for development are internally determined by the content and structure of the technology as a body of knowledge. How new designs are introduced, how experience is gained and utilized, and how information is obtained, are the issues more relevant to directing the path that early technological development takes, once a project has been initiated. The later implementation and adoption stages of the technology are the periods that are being more strongly influenced by the wider socio-economic and political context.

## 6.3 Important elements of the process

### 6.31 Idea generation

The first stages of the process of technological innovation involves the conception and evaluation of ideas. Often this is characterized by a loose informal organization with a minimum of constraints or control. The generation of an idea whether it concerns a new product, process or system, is mainly a problem solving process of a psychological nature. The problem, in general, is a recognition of discrepancies between a desired situation and the existing one. This must be of sufficient magnitude to be worthy of consideration. The idea generation stage can thus be modeled as the fusion between perception of needs (the problem), and the means to fulfil them.

Both the definition of the problem and potential solutions depend on information. In practice, attention is most often focussed on the means or problem solving aspects than on the problem definition. The character of the technological possibilities that emerge depend on both the kind of information that is received and where the emphasis may be placed. The information utilized may be in the form of knowledge received through studies, practical experience etc., or information collected through printed material or discussions with knowledgeable individuals (we discuss the implication of the source and flow of information later in this paper).

The definition of the problem itself also depends on information about what is needed. Neglecting the provision of information about needs, in any technological study may easily give rise to a loosely defined problem. The consequence may be unfortunate since the outcome of the idea generation stage has a decisive influence on the innovation process and the path eventually taken.

In our protocol study, the problem was identified as the need to harness the vast communication potential of optical networks, brought about by the technological developments in areas wide high capacity optical fibers and narrow-bandwidth lasers. By understanding the salient characteristics of

the optical communication environment and the approaches adopted to deal with such an environment in the satellite communication arena, we established the basic means to effect the desired result.

## 6.32 Idea realization

The decision to go ahead with an idea or proposal represents a key point in the innovative process. The idea generation process is brought to an end; the utility of further research is established and the ideas have been converted into a project. This represents a significant change with regards to activities, focus and organization and structure of the undertaking. Whereas idea generation can be characterized as software, the realization is concerned with hardware. The idea for a solution must now be converted into a working reality. This requires a more structured organization and an orderly procedure with a greater amount of formal exploration of issues and process iterations.

Within our study, the idea realization stage involved a structured approach to understanding the nature of satellite communication. This involved interpreting and comparing the many proposed schemes which sought to increase the effectiveness of the multiaccess communication. This was followed by a comparative evaluation of the performances of these schemes.

## 6.33 Analytical problem solving

A common model of problem solving involves the steps: 1) definition of the problem, 2) collection of data, 3) analysis of data, 4) development of alternatives, and 5) selection and implementation of solutions. As mentioned, the first step is very important in determining the scope and direction of the steps that follow. Defining the problem must take into consideration the needs, constraints, and possible objective priorities.

In our study, the data collection involved an exhaustive survey of satellite protocol proposals and this was followed by the analysis of the operation and performance of the various methods. By establishing the most successful schemes based on certain defined criteria, the alternatives for optical protocol

methods was developed. We were then able to select certain approaches and proceed to analyze in closer detail their suitability in achieving our overall objective. This process of examining possible methods resulted in a number of unworkable solutions but with each a growing understanding of the nature of the problems that had to be solved. This iterative process of analysis, learning, modifying assumptions and refocussing emphases, proved to be most instructive in pointing the direction towards the eventual protocol solution.

## 6.4 Factors influencing Technological Innovation

### · 6.41 The Processes of Learning and Scaling

Sahal [Sahal 81] in his analysis of technological innovations points to two key determinants of technical progress. The first is the process of learning or the acquisition of knowledge and understanding through experience. The other is the process of scaling or the patterning of systems to perform certain desired tasks. These processes are seen to play a pivotal role in the development of technologies. Indeed, technological development is dependent on the interplay of myriad factors. However it is felt that these main processes of learning and scaling affect a wide variety of the determining variables. As we examine more closely the area of optical communication protocol design, we can begin to understand the influence that these factors have had on the development of the technology.

The process of learning goes beyond the mere acquiring of knowledge of system principles and theories. Rather it embraces the more important aspects of understanding through experience and experimentation. Technical progress depends not only on learning from from past failures but also on learning to anticipate new opportunities. It is not only a matter of determined effort in an attempt to operate within system limitations, but also a matter of abandoning a chosen approach once there is a clear indication of its deficiencies. It involves a constant striving to attain the possibilities that are intrinsic to the system.

Within the first stage of our work with satellite system protocols we were able to gain better understanding of the important aspects of protocol design in high propagation delay environments. The methods applied to solving the problems of access coordination. This knowledge of underlying constraints was an appropriate first step, the greatest amount of learning however only came with the early attempts to devise new schemes. While it was relatively easier to propose operating protocols once we began to apply some of these, the difficulties became clearer. As more approaches proved unsuccessful due to theoretical or practical system constraints, the more the experience gained pointed towards certain classes of schemes. The element of trial and error entered in a fundamental way in the development of new techniques. This is not to imply an ad hoc approach, but rather to highlight the fact that experimentation was an inescapable part of the design process. The experience gained from each design attempt served to better define the framework in which schemes would be most successful.

The focal point in these learning processes has been changes in the scale of the object system. While conventional wisdom has it that the scale of an object depends on the availability of technology, this is a somewhat over simplified view. Indeed, barriers to economies of scale often prove to be temporary with the advent of relevant innovations. However, it is seldom possible to change the scale of a system and not affect its form and structure. Thus scaling the operating performance of a system is an important design problem in itself. Scaling is therefore not simply a by-product of successful innovation but rather innovation often originates during the course of successful changes in scale. As well technological innovation may result when attempts at scale changes encounter the intrinsic limitations of existing technologies.

The element of scaling has been a principal driver in the quest to design more effective optical communication protocols. The potential ability of optical fibers and emerging laser technologies to facilitate the transports vast quantities of information can only be harnessed through the application of appropriate communication access methods. The order of magnitude increase in this communication

capacity means that the previous single channel technologies are no longer as effective on a system level. The increased communication speeds/bandwidth has posed new problems. The changes in scale has resulted in a new multichannel environment in which the experience acquired in the design of protocols for satellite systems and local area networks is not wholly transferable.

## 6.42 Evolution rather than Revolution

One of the established concepts in the study of technological development is the evolutionary nature of the process. This is based on the principle that successful technology tends to create an operating paradigm of its own so influencing future development. The process thus becomes self-generating as well as self-constraining. This strong influence can be traced to both the ability of a successful technology to satisfy its present demands, and the constituent of followers that it produces. As a result, technical progress often moves in bit-by-bit modifications rather than by the development of alternative techniques. Newer technologies are thus often the cumulative sum of the gradual refinements of an older technology. As Kuhn [Kuhn, 70 p.163] explains it:

> ... once the reception of a common paradigm has freed the scientific community from the need constantly to re-examine its first principles, the members of that community can concentrate exclusively upon the subtlest and most esoteric of the phenomena that concern it.

That is not to deny the occurrence of radical or significantly altering developments, but to point out that the frequency of occurrence of such innovations is often overstated. There is a general tendency to overstudy exciting, radical innovations. This was also concluded in an analysis by Clark and Staunton [Clark and Staunton, 89] which suggested that many innovations that are essentially entrenching on-going directions, have been wrongly reported as radical.

In our study of satellite system protocols we have been able to trace such incremental developments in a number of the schemes proposed. A few key innovations have been the source of the many approaches adopted. Often, as we have described in the survey of chapter 3, the methods represent simple variations on a central theme. The field of access protocols has moved from fixed assignment

approaches to demand assignment schemes. The success of random access methods initiated with the ALOHA protocol and the use of reserved access have been the foundations from which all later schemes have been derived. Even as we enter the area of optical communication schemes we can clearly see the influence of schemes based on the understood principles of satellite communication. It was in fact such a premise of better understanding of established single channel communication protocols, that drove our own new design.

As we now examine that process which was based on sound technical reasoning, it is recognized that such an approach may have pre-constrained the results that could be derived. This reliance on proven concepts from the past as well as the ability of some techniques to be adapted to their environments makes it more likely that such techniques can be a vehicle for further advances. As Sahal [Sahal, 81 p.310] points out:

> Initially, the focal point of innovative activity tends to be the adaptation of technology to the task environment; subsequently, it is the adaptation of the task environment to the technology.

Technical change is thus often characterized as a cumulative process of minor improvements; incremental innovations with only occasional major or radical innovations. These radical innovations serve as the base of a technological trajectory. These trajectories of development dictate the progression of technology often perpetuating technological evolution whereby technology builds upon technology in a step-by-step manner.

The study of the satellite protocol literature show the above findings of technological incrementalism to be particularly true. The foundation was laid with the introduction of the ALOHA scheme and continues with approaches that modify the basic operation. In the case of optical communication protocols the early work has shown the framing of the problem in terms of the established single channel communication concepts. While this is in part due to the constraints of electronic communication, it is interesting to note that no new, radically different approach to this more complex

problem has as yet surfaced. Understanding the foundations and success of relevant technologies thus provides useful insight into how future applications may evolve.

### 6.43 Sources of Information and Information Flow

In studying the research and development process some analysts have focussed on the influence of idea and information flow [Kelly and Kranzberg, 78]. They have traced the flow of technical information and its impact on the creation and development of ideas. One set of findings that emerge from this idea flow literature is the observed differences between the information seeking practices of basic and applied researchers; reflecting differences about objectives and time horizons. Basic researchers spend more time formulating and defining problems and their information sources tend to be peer-review journals and conference proceedings. The applied researcher by contrast, works with a well defined problem and must quickly find an acceptable solution. The applied researcher will frequently communicate with fellow project team members and other experts within his/her organization. Use of printed material is often limited to in-house technical reports [Marquis and Allen, 67].

Utterback [Utterback, 71] also highlighted in his work the influence of information on the process of technical innovation. In his stage process model of innovation, idea generation represents the initial phase in the development of new technology. The origination of ideas for new products occur in response to recognition of a need or problem or, less frequently, in response to recognition of an existing technical possibility or means. This has also been defined in some of the innovation studies literature as market pull versus technology push (where market pull implies a more direct recognition of the economic possibilities for a given technology). A major finding of his study was that two pieces of information were associated with each idea. The recognition of a need, problem or opportunity, and recognition of a means or technique by which to satisfy the need, solve the problem, or meet the opportunity.

Further hypotheses were supported by his work. Firstly, the technology that is used in innovations stimulated by technical possibilities is likely to more recent in origin than technology employed in

need-stimulated innovations. This is based on the implication that older technology is less likely to attract attention. New discoveries on the other hand may receive wider communication and attention and as a result provide a focus for innovative activity.

The effects of these reported factors can be seen from our search for a protocol for high speed optical networks. The project was based on a recognized need given the dearth of proposed methods to deal with the problem of multiple channel wavelength division optical communication. The study thus arose out of an understood need for an appropriate technical solution. Without the stimulus of new technology creating an opportunity, the problem solving stage of the process relied on the use of older established technologies of satellite communication;seen as a technique that could be applied in satisfying the particular need.

In addition, the project being more basic research in nature also supported the conclusion that sources of information employed in the problem solving stage of the process would comprise principally of peer-review journals and conference proceedings. For our part this source of information was seen as providing up-to-date records of the theoretical studies conducted in the field. The nature of academic institutions also only provided limited forum for exposure to more informal sources of information. Through seminars, however, and discussions with those more exposed to the field beyond the university, it was possible to gain a better practical understanding of some of the difficulties to be resolved in the protocol design problem. It was clearly recognized that by focussing on journals and conference reviews, many of the issues of practical system implementation and practical design were not readily appreciated.

## 6.44 The academic organization

The role of organizations in the innovation process is one that has received considerable study. These have focussed on the dominant general theories of organizations [Perrow, 79]. This extends from the classical approach that premises the clear goals of organizations, to the human relations model which highlights informal interactions within organizations, to the contingency theory approach that

emphasizes the adaptability of the organizational structure to the tasks to be undertaken, to the systems theory approach that looks at organizations as goal-oriented systems with elements made up of individuals and their work apparatus.

The academic institute is indeed a goal-oriented organization of which much analysis can be done. The effects of the structure on the development of innovation at the level of the individual student may, however, be very subtle beyond the phase of problem selection. Of more significance is the role played by academic advisors and research support staff. In this regard professors have dominating influences in their positions as technical gatekeepers and technical entrepreneurs in the classical sense. The gatekeeper by his/her ability to identify scientific and technical information of relevance to the activities of the student researcher is able to profoundly affect the course of innovation. By championing a particular activity, the technical entrepreneur also motivates a particular direction for technological development.

Competition among individuals within organizations has often been sighted as a important source of generating innovative activity. The influence of the competitive pressures among individuals within the academic organization, however, may not be generally felt due to the often independent nature of the undertakings. Also, the fact that there are no reward structures designed to foster overt competition among individuals. The organizational environment is nonetheless one that encourages achievement both through peer achievement as well as the opportunity to combine creativity with more routine learning. The general diversity of the organization and the understood goal of fostering new ideas are factors that play a subtle role in the innovation process.

## 6.5 Implications and recommendations

As discussed, scale in the innovation process is not at the peripheral but rather it is often the crux of research and development at the grass roots level. A great deal of learning in research and development activity has to do with solutions to problems posed by contemplated changes in scale of systems or technologies. Scale and learning are thus important determinants of innovation and are both intertwined elements of the process. Changing the scale of technology may do not always involve the simple manipulation or extrapolation of system elements. An alteration of the technology may be required to deal with the new conditions encountered. In such an environment innovation may emerge.

· An integral part of the innovation process is the experience gained through direct experimentation and analysis. The least successful alternatives may often produce the best understanding of the constraints involved in the problem being addressed, or in the solution targeted. There can thus be no substitution for a broad consideration of alternatives in the problem solving stage of the process to allow for many avenues to be pursue and valuable experience to be gained.

The understanding that new technological innovation evolves from established technologies is important in structuring new technical design studies. The recognition of biases in any process is the first step to arriving at unconstrained results. Understanding the importance of certain technological developments provides a means to forecast later related technology developments. In understanding the evolution of technologies that results from a major technological innovation, we have a perspective that can allow us to look beyond the type of evolutionary changes that usually follows. This leads to much broader thinking that can provide significantly different developments.

In this chapter we have pointed to the importance of information flow from external sources in the developing of practical systems. There is a recognized need to expose the research community to sources of technical information emerging from outside the organization especially with regards to industrial

practice. Basic researchers can learn to better appreciate practical considerations in their designs through exposure to the work and understanding gained by applied researchers and professionals. It is important therefore to the academic research community that opportunities are created for formal and informal interaction with others involved, not only in research, but also in design and development.

# Chapter 7

# Summary and Conclusions

In our study of satellite protocol schemes we were able to identify certain basic characteristics and trends in the many approaches proposed. At one extreme of protocol operation was the random access schemes which offered the advantages of easy implementation without the need for coordination overhead. At the other end of the spectrum was the fully coordinated reservation approaches in which channel access was tightly controlled. However, while removing the possibility of collisions, the need to make reservations and await assignment increased the average delay times of these schemes. In addition, there is a greater overhead requirement in the maintenance of system information. Many of the earlier approaches exhibited characteristics that could be classified as predominantly random access or reservation access.

A more efficient approach to the multiple access problem came with the combining or integration of the two extreme classes. These schemes were able to take advantage of the short delay times of random access approaches during periods of low system load yet move to more of a reservation base as the system became more heavily loaded. This integration was mainly achieved by assigning random access and reserved access periods within a transmission frame or through the use of asynchronous reservation periods which were adaptive to system conditions. The latest means to achieving this integration has been the transmission of reservation information during normal random access transmission to be used in the event of collisions. These schemes while displaying good throughput-delay performances are

nonetheless more complex in their implementation. The performance results compared clearly showed the superiority of integrated access schemes.

Our early attempts to incorporate the principles of mixed protocol operation in the design of a new optical scheme, brought more clearly to light the additional complexity involved in the multichannel environment. Given realistic hardware constraints, the employing of coordinated access schemes could not be effectively accommodated. Through a number of failed attempts to configure a suitable coordinated scheme the experience gained from these attempts pointed to the use of a random access based approach. The XURN scheme was the result of this new direction, providing a scheme that was random access in basis, yet capable of high throughput and a flexible adaptive operation.

The analysis of this proposed scheme confirmed it potential for multichannel optical network operation. For the 10 station network case analyzed, the scheme is able to provide a maximum throughput of 0.7 which is significantly better than that obtained from the analysis of other optical protocols for similar sized networks. The scheme requires a calculation overhead to operate effectively, but this factor is not expected to present an insoluble implementation problem.

# Appendix A

In deriving the transition probability matrix for the analysis of the system operation we begin by calculating the probability of successful channel transmissions.

From the expression 5.41.06 (p. 107) we are able to calculate $f(n, q)$, the probability of $q$ successful channel transmissions given $n$ ready stations in the network. Given the number of network stations and available channel frequencies, these probabilities can be exactly calculated inspite of the absence of an explicit closed form expression. The results of this calculation for values of $n$ are thus given below:

Since $b = 4$,      $f(n, q) = 0$      for all $q > 4$,
and in general,   $f(n, q) = 0$      for all $q > n$.

Hence,

| For   $n = 0$ | | For   $n = 1$ | For   $n = 2$ |
|---|---|---|---|
| $f(0,0) = 1$ | | $f(1,0) = 1/5$ | $f(2,0) = 5/45$ |
| $f(0,q) = 0$ | for all $q > 0$ | $f(1,1) = 4/5$ | $f(2,1) = 16/45$ |
| | | | $f(2,2) = 24/45$ |

| For   $n = 3$ | For   $n = 4$ | For   $n = 5$ | For   $n = 6$ |
|---|---|---|---|
| $f(3,0) = 1/15$ | $f(4,0) = 5/105$ | $f(5,0) = 3/63$ | $f(6,0) = 1/210$ |
| $f(3,1) = 4/15$ | $f(4,1) = 24/105$ | $f(5,1) = 12/63$ | $f(6,1) = 24/210$ |
| $f(3,2) = 6/15$ | $f(4,2) = 36/105$ | $f(5,2) = 24/63$ | $f(6,2) = 90/210$ |
| $f(3,3) = 4/15$ | $f(4,3) = 32/105$ | $f(5,3) = 16/63$ | $f(6,3) = 80/210$ |
| | $f(4,4) = 8/105$ | $f(5,4) = 8/63$ | $f(6,4) = 15/210$ |

| For   $n = 7$ | For   $n = 8$ | For   $n = 9$ | For   $n = 10$ |
|---|---|---|---|
| $f(7,0) = 0$ | $f(8,0) = 0$ | $f(9,0) = 0$ | $f(10,0) = 0$ |
| $f(7,1) = 1/30$ | $f(8,1) = 0$ | $f(9,1) = 0$ | $f(10,1) = 0$ |
| $f(7,2) = 9/30$ | $f(8,2) = 6/45$ | $f(9,2) = 0$ | $f(10,2) = 0$ |
| $f(7,3) = 15/30$ | $f(8,3) = 24/45$ | $f(9,3) = 2/5$ | $f(10,3) = 0$ |
| $f(7,4) = 5/30$ | $f(8,4) = 15/45$ | $f(9,4) = 3/5$ | $f(10,4) = 1$ |

In equation 5.41.09 (p. 113) we defined $g(q, s)$ as the probability of exactly $s$ successful packet receptions given $q$ successful channel transmissions. Using that equation, we obtain the following results:

$g(0,0) = 1,$

$g(q,s) = 0$  for all $s > q$

and  $g(q,0) = 0$  for all $q > 0$,  since one packet is always correctly received for multiple packets arriving at any station.

additionally,

$g(1,1) = 1$

$g(2,1) = 8/81$  $g(3,1) = 7/729$  $g(4,1) = 2/2187$
$g(2,2) = 73/81$  $g(3,2) = 195/729$  $g(4,2) = 134/2187$
  $g(3,3) = 527/729$  $g(4,3) = 940/2187$
    $g(4,4) = 1111/2187$

Consistency is confirmed by ensuring that the probabilities in each of these columns like those above always sum to 1.

These successful reception probabilities must now be factored together with the probability of successful channel transmissions using equation 5.41.10 (p. 114) in order to obtain the probability $P_t(n, s)$ of exactly $q$ completely successful packet transmissions in a given slot, for a specified number of busy stations $n$. The following results are obtained:

$P_t(0,0) = 1$

$P_t(1,0) = 1/5$
$P_t(1,1) = 4/5$

$P_t(2,0) = 135/1215$
$P_t(2,1) = 496/1215$
$P_t(2,2) = 584/1215$

$P_t(3,0) = 0.0667$
$P_t(3,1) = 0.3087$
$P_t(3,2) = 0.4318$
$P_t(3,3) = 0.1928$

$P_t(4,0) = 0.0476$
$P_t(4,1) = 0.2654$
$P_t(4,2) = 0.3952$
$P_t(4,3) = 0.2531$
$P_t(4,4) = 0.0387$

$P_t(5,0) = 0.0476$
$P_t(5,1) = 0.2307$
$P_t(5,2) = 0.4190$
$P_t(5,3) = 0.2382$
$P_t(5,4) = 0.0645$

$P_t(6,0) = 0.0048$
$P_t(6,1) = 0.1603$
$P_t(6,2) = 0.4925$
$P_t(6,3) = 0.3061$
$P_t(6,4) = 0.0363$

$P_t(7,0) = 0$
$P_t(7,1) = 0.0679$
$P_t(7,2) = 0.4143$
$P_t(7,3) = 0.4331$
$P_t(7,4) = 0.0847$

$P_t(8,0) = 0$
$P_t(8,1) = 0.0186$
$P_t(8,2) = 0.2832$
$P_t(8,3) = 0.5288$
$P_t(8,4) = 0.1693$

$P_t(9,0) = 0$
$P_t(9,1) = 0.0044$
$P_t(9,2) = 0.1438$
$P_t(9,3) = 0.5471$
$P_t(9,4) = 0.3048$

$P_t(10,0) = 0$
$P_t(10,1) = 2/2187$
$P_t(10,2) = 134/2187$
$P_t(10,3) = 940/2187$
$P_t(10,4) = 1111/2187$

Consistency checks are easily applied as the sum of each of these prob. columns must be 1.

Our next consideration is the evaluation of packet arrivals at network stations. From equation (5.41.12) we are able to derive the probability of $j$ packet arrivals at non-ready stations for different values of the arrival rate at each station. The packet arrival and success probabilities can then be combined using equation 5.41.13 (p. 115), to derive exact values for the elements of the transition probability matrix for different packet arrival rate values. Again a consistency check must confirm that the row of the transition probability matrix all sum to 1. The results were obtained by programming the appropriate equation in *Mathematica* software. The following is an example of the particular results obtained for the arrival rate at each station, $\lambda/M = 0.2$. Rounding off results the (11x11) transition probability matrix, which defines state 0 to state $M$ (=10) operation, is:

{0.1353, 0.2996, 0.2985, 0.1762, 0.0683, 0.0181, 0.0004, 0, 0, 0, 0}

{0.1322, 0.2966, 0.2992, 0.1789, 0.0702, 0.0189, 0.0035, 0.0005, 0, 0, 0}

{0.0971, 0.2543, 0.3016, 0.2118, 0.0972, 0.0304, 0.0065, 0.0010, 0, 0, 0}

{0.0475, 0.1802, 0.2901, 0.2621, 0.1483, 0.0553, 0.0023, 0.0002, 0, 0, 0}

{0.0117, 0.0917, 0.2289, 0.2966, 0.2250, 0.1064, 0.0324, 0.0064, 0.0008, 0, 0}

{0, 0.0237, 0.1139, 0.2628, 0.3010, 0.1968, 0.0788, 0.0197, 0.0030, 0.0003, 0}

{0, 0, 0.0163, 0.1520, 0.3479, 0.3092, 0.1370, 0.0330, 0.0043, 0.0003, 0}

{0, 0, 0, 0.0465, 0.2686, 0.3921, 0.2238, 0.0608, 0.0079, 0.0004, 0}

{0, 0, 0, 0, 0.1135, 0.4047, 0.3524, 0.1139, 0.0148, 0.0006, 0}

{0, 0, 0, 0, 0, 0.2496, 0.5032, 0.2170, 0.0297, 0.0008, 0}

We can see that the sum of any of the rows of this matrix [P] is 1. The matrix also satisfies the expected form described on page 116 (lower left hand corner zeros based to the maximum number of channel successes that can occur in a given slot period). The steady state probabilities of the Markov chain can be obtained by solving the equation:     $\pi = \pi [P]$

or by taking any of the row of $[P]^n$ for very large n, since the power of $[P]$ converges to a matrix of equal rows which defines $\pi$.

For this case example we obtained (again rounding off the results for presentation purposes)

$\pi$ = {0.0765, 0.2081, 0.2705, 0.2263, 0.1358, 0.0590, 0.0186, 0.0039, 0.0005, 0, 0}

Here again the accuracy of the results was be confirmed by the elements of the row vector summing to 1.

Applying these steady state probabilities to the throughput and delay equation 5.41.16 and 5.41.17 respectively (p. 117), we obtain a throughput delay point {2.3241, 1.4686}.

· Other points are derived by the same process for other values of the packet arrival rate at each network station.

# References

[Abr 70]     N. Abramson, "The ALOHA system-Another alternative for computer
             communications," in *AFIPS Conf. Proc.*, 1970 *Fall Joint Comput.* vol. 37, pp. 281-285.

[Abr 73]     N. Abramson, "Packet Switching with Satellites," in *AFIPS Conf. Proc.*, vol. 42, pp.
             695-702, June 1973.

[Aca 87]     A. Acampora, "A multi-channel multihop local lightwave network," *Proc.*
             *GLOBECOM 1987*, pp. 37.5.1-37.5.9, Tokyo, Japan 1987.

[BF 78]      F. Borgonovo and L. Fratta, "SRUC: A Technique for Packet Transmission on Multiple
             Access Channels," in *Proc. Intl. Conf. Comput. Commun.*, Kyoto, Japan, Sept. 1978, pp.
             601-607.

[Bin 75]     R. Binder, "A dynamic packet switching system for satellite broadcast channels," in
             *Proc. ICC'75*, San Francisco, CA, June 1975.

[BP 79]      M. Balagangadhar and R. Pickholtz, "Analysis of a Reservation Multiple Access
             Technique for Data Transmission via Satellites," *IEEE Trans. Commun.*, vol. COM-27,
             no.10, pp. 1467-1475, October 1979.

[Bra 88]     C. Brackett, "Dense WDM Networks," in *Conf. Proc., IEE ECOC '88*, Brighton, UK, 1988.

[Cap 79]     J. I. Capetanakis, "Tree algorithms for packet broadcast channels," *IEEE Trans. Inform.*
             *Theory*, vol. IT-25, pp. 505-515, Sept. 1979.

[CBK 89]     D.M. Chitre, A.C. Briancon, and R. Kohli, "Random Access with Notification - A New
             Multi-Access Scheme for VSAT networks," *COMSAT Tech. Rev.*, vol. 19, no. 1, pp. 99-
             121, Spring 1989.

[CDR 90]     M-S Chen, N. Dono and R. Ramaswami, "A Media Protocol for Packet Switched
             Wavelength Division Multiple-access Metropolitan Area Networks", *IEEE J. Select.*
             *Areas Commun.*, vol. 8, no. 6, pp. 1048-1057, August 1990.

[CG 88a]    I. Chlamtac, A. Ganz, and Z. Koren, "Prioritized Demand Assignment Protocols and Their Evaluation," *IEEE Trans. Commun.*, vol. COM-36, no. 2, pp. 133-143, Feb. 1988.

[CG 88b]    I. Chlamtac and Aura Ganz, "Design and Analysis of Very High Speed Network Architectures," *IEEE Trans. Commun.*, vol. COM-36, no.3, pp. 252-262, March 1988.

[CG 88c]    I. Chlamtac and Aura Ganz, "Channel Allocation Protocols in Frequency-Time Controlled High Speed Networks," *IEEE Trans. Commun.*, vol. COM-36, no. 4, pp. 430-440, April 1988.

[CG 88d]    I. Chlamtac and Aura Ganz, "A Multibus Train Communication (AMTRAC) Architecture for High-Speed Fiber Optic Networks," *IEEE J. Select. Areas Commun.*, vol. 6, no. 6, pp. 903-912, July 1988.

[CG 89-1]    I. Chamtac and A Ganz, "Design Alternatives of Asynchronous WDM Star Networks", Univ. of Massachusetts, 1989.

[CG 89-2]    I. Chamtac and A Ganz, "Path Allocation Access Control in Fiber Optic Communication Systems", *IEEE Trans. Comput.*, vol. 38, no. 10, pp. 1372-1381, October 1989.

[CN 78]    W. Chu and W. Naylor, "Measurement and Simulation results of C-PODA protocol performance," NTC 78 Conf. Rec., Nastional Telecommunications Conference, Birmingham, Alabama, December 1978.

[CR 83]    G. Choudhury and S. Rappaport, "Diversity ALOHA - A Random Access Scheme for Satellite Communications," *IEEE Trans. Commun.*, vol. COM-31, no. 3, pp. 450-457, March 1983.

[CRW 73]    W.R Crowther, R. Rettberg, D. Walden, S. Ornstein and F. Heart, "A system for broadcast communication: Reservation-ALOHA," in *Proc. 6th Hawaii Int. Syst. Sci. Conf.*, January 1973.

[GC 87]    A. Ganz and I. Chlamtac, "Finite Buffer Queueing Models for Single and Multiple Transmissions Channel Access Protocols," *Proc.* 2nd International MCPR Workshop, Rome, Italy, May 1987, pp. 321-335.

[GE 81]    E. P. Greene and A. Ephremides, "Distributed Reservation Control Protocols for Random Access Broadcasting Channels," *IEEE Trans. Commun.*, vol. Com-29, no. 5, pp. 726-735, May, 1981.

[GP-K 82]    L. Georgiadis and P. Papantoni-Kazakos, "A Collision Resolution Protocol for Random Access Channels with Energy Detectors," *IEEE Trans. Commun.*, vol. COM-30, no. 11, pp. 2413-2420, November 1982.

[GSS 82]    D. Guha, D. Schilling, and T. Saadawi, "Dynamic Reservation Multiple Access Technique for Data Transmission via Satellites," in *Proc. IEEE Infocom 82*, pp. 53-61.

[HKS 87]    I.M Habbab, M. Kavehrad, and C-E.W Sundberg, "Protocols for Very High-Speed Optical Fiber Local Area Networks Using a Passive Star Topology," *J. Lightwave Technol.*, vol. LT-5, no. 12, pp.1782-1793, December 1987.

[Hum 90]    P.A Humblet, "Protocols for High Speed Wavelength Division Optical Networks-Proposal", MIT, January 1990.

[JB 78]    I.M. Jacobs, R. Binder, and E.V. Hoversten, "General Purpose Packet Satellite Networks," *Proc. IEEE*, vol. 66, no. 11, pp. 1448-1467,

[KH]    R.s Kennedy and P.A Humblet, "Open Questions concerning Network Design in an era of Excess Bandwidth", MIT.

[KY 78]    L. Kleinrock and Y. Yemini, "An optimal adaptive scheme for multiple acess broadcast communication," in *Proc. Int. Conf. Commun.*, 1978, pp. 7.2.1-7.2.5.

[Lam 77]    S. Lam, "Satellite multi-access schemes for data traffic," in *Proc. Int. Conf. Commun.*, Chicago, IL, 1977, pp37.1-19 - 37.1-24

[Lam 79]    S. Lam, "Satellite Packet Communication - Multiple Acccess Protocols and Performance," *IEEE Trans. Commun.*, vol. Com-27, no.10, pp. 1456-1466, October 1979.

[LK 75]    S. Lam and L. Kleinrock, "Dynamic control schemes for a packet switched multi-access broadcast channel," in *AFIPS Conf. Proc.*, vol. 44, AFIPS Press, Montvale, N.J. 1975, pp. 143-153.

[LM 83]  H. Lee and J. Mark, "Combined Random/Reservation Access for Packet Switched Transmission Over a Satellite with On-Board Processing; Part I - Global Beam Satellite," *IEEE Trans. Commun.*, vol. Com-31, no. 10, pp. 1161-1171, October 1983.

[LT 83]  T. T. Liu and D. Towsley, "Window and Tree Protocols for Satellite Channels," in *Proc. IEEE Infocom, 83*, pp. 215-221.

[Mar 78]  J. Mark, "Global Scheduling Approach to Conflict-Free Multiaccess via a Data Bus," *IEEE Trans. Commun.*, vol. Com-26, no. 9, pp. 1342-1351, September, 1978.

[Meh 90]  N. Mehravari, "Performance and Protocol Improvements for Very High Speed Optical Fiber Local Area Networks Using a Passive Star Topology," *J. Lightwave Technol.*, vol. 8, no. 12, pp. 520-530, April 1990.

[MK 83]  L. Merakos and D. Kazakos, "Multiaccess of a slotted channel using a control mini-slot," in *Proc Int. Conf. Commun.*, 1983, pp. C5.3.1-C5.3.6.

[NM 77]  S.F. Ng and J. W. Mark, "A multiaccess model for packet switching with a satellite having some processing capability," *IEEE Trans. Commun.*, vol. Com-25, no. 1, pp. 128-135, January, 1977.

[Ray 84]  D. Raychaudhuri, "ALOHA with Multipacket Messages and ARQ-Type Retransmission Protocols - Throughput Analysis," *IEEE Trans. Commun.*, vol. COM-32, no. 2, pp. 148-154, February 1984.

[Ray 85]  D. Raychaudhuri, "Announced Retransmission Random Access Protocols," *IEEE Trans. Commun.*, vol. COM-33, no. 11, pp. 1183-1190, November 1985.

[Ray 87]  D. Raychaudhuri, "Stability, Throughput, and Delay of Asynchronous Selective Reject ALOHA," *IEEE Trans. Commun.*, vol. COM-35, no. 7, pp. 767-772, July 1987.

[Riv 85]  R. Rivest, "Network Control by Bayessian Broadcast," Report MIT/LCS/TM-285. Cambridge, MA: MIT, Laboratory for Computer Science.

[Rob 73]  L. Roberts, "Dynamic allocation of satellite capacity through packet reservation," in *AFIPS Conf. Proc.*, vol. 42, June 1973.

[Rub 77a]     I. Rubin, "Integrated Random-Access Reservation schemes for multi-access communication channels," School Eng. Appl. Sci., Univ. California, Los Angeles, Tech. Rep. UCLA-ENG-7752, July 1977.

[Rub 77b]     I. Rubin, "A Group Random-Access Procedure for Multi-Access Communication Channels," in NTC'77 Conf. Rec. Nat. Telecommun. Conf., Los Angeles, CA, Dec. 1977, pp. 12:5-1 - 12:5-7.

[Rub 79]      I. Rubin, "Access-Control Disciplines for Multi-Access Communication Channels: Reservation and TDMA Schemes," IEEE Trans. Inform. Theory, vol. IT-25, no. 5, pp. 516-536, Sept. 1979.

[SC 88]       F. Scholl and M. Coden, "Passive Optical Star Systems for Fiber Optic Local Area Networks," IEEE J. Select. Areas Commun., vol. 6, no. 6, pp. 913-923, July 1988.

[SKP 86]      E. Sykas, D. Karvelas, and E. Protonotarios, "Queueing Analysis of Some Buffered Random Multiple Access Schemes," IEEE Trans. Commun., vol. Com-34, no. 8, pp. 790-798, August 1986.

[Tas 84]      S. Tasaka, "Multiple-Access Protocols for Satellite Packet Communication Networks: A Performance Comparison," Proc. IEEE, vol. 72, no. 11, pp. 1573-1582, November 1984.

[TC 86]       D. Tsai and J-F. Chang, "Performance Study of an Adaptive Reservation Multiple Access Technique for Data Transmsission," IEEE Trans. Commun., vol. Com-34, no. 7, pp. 725-727, July 1986.

[TI 86]       S. Tasaka and K. Ishida, "The SRUC Protocol for Satellite Packet Communication - A Performance Analysis," IEEE Trans. Commun., vol. Com-34, no. 9, pp. 937-945, September 1986.

[TK 78]       F. Tobagi and L. Kleinrock, "Packet Switching in Radio Channels: Part III - Polling and (Dynamic) Split-Channel Reservation Multiple Access," IEEE Trans. Commun., vol. Com-24, no. 8, pp. 832-845, August, 1976.

[Tob 80]      F.A Tobagi, "Multiaccess Protocols in Packet Communication Systems", IEEE Trans. Commun., vol. COM-28, pp. 468-488, April 1980.

[WE 80]  J.E Wieselthier and A. Ephremides, "A New Class of Protocols for Multiple Access in Satellite Networks", *IEEE Trans. Automat. Contr.*, vol. AC-25, pp. 865-879, Oct. 1980.

[WGT 86]  C.J. Wolejsza, D. Taylor, M. Grossman, and W.P. Osborne, "Multiple Access Protocols for Data Communications via VSAT Networks," *IEEE Commun. Mag.*, vol. 25, no. 7, pp. 468-488, July 1986.

[Yum 87]  T.P. Yum, "The Design and Analysis of the Scheduled-Retransmission Multiaccess Protocol for Packet Satellite Communications," *IEEE International Conference on Commun.*, Seattle, Washington, 1987, Proc., pp. 278-283, November 1978.

· [Bertsekas and Gallager, 87]  D. Bertsekas and R. Gallager, *Data Networks*, Prentice-Hall, Inc. New Jersey, 1987.

[Clark and Staunton, 89]  P. Clark and N. Staunton, *Innovation in Technology and Organization*, Routledge, London, 1989.

[Holt, 87]  K. Holt, *Innovation: A Challenge to the Engineer*, Elsevier, Amsterdam, 1987.

[Johnston and Gibbons, 75]  R. Johnston and M. Gibbons, "Characteristics of Information Usage in Technological Innovation, *IEEE Trans. Eng. Mgmt.*, vol. EM-18, no. 4, Nov. 1971.

[Kelly and Kranzberg, 78]  P. Kelly and M. Kranzberg, *Technological Innovation: A Critical Review of Current Knowledge*, San Francisco University Press, 1978.

[Kuhn, 70]  T. Kuhn, *The Structure of Scientific Revolutions*, The University of Chicago Press, 1970.

[Marquis and Allen, 67]  D. Marquis and T. Allen, "Communication Patterns in Applied Technology," *AmericanPsychologist*, 21:1052-1060, 1967.

[Mole and Elliot, 87]  V. Mole and D. Elliot, *Enterprising Innovation: An Alternative Approach*, Frances Pinter (Publishers), London, 1987.

[Sahal, 81]              D. Sahal, _Patterns of Technological Innovation_, Addison-Wesley, Massachusetts, 1981.

[Tornatzky, Eveland _et al_, 83]  L. Tornatzky, J. Eveland, M. Boylan _et al_, _The Process of Technological Innovation: Reviewing the Literature_, Productivity Improvement Research Section, Division of Industrial and Technological Innovation, National Science Foundation, May 1983.

[Utterback, 71]         J. Utterback, "The Process of Innovation: A Study of the Origination and Development of Ideas for New Scientific Instruments," _IEEE Trans. Eng. Mgmt._, vol. EM-18, no. 4, Nov. 1971.