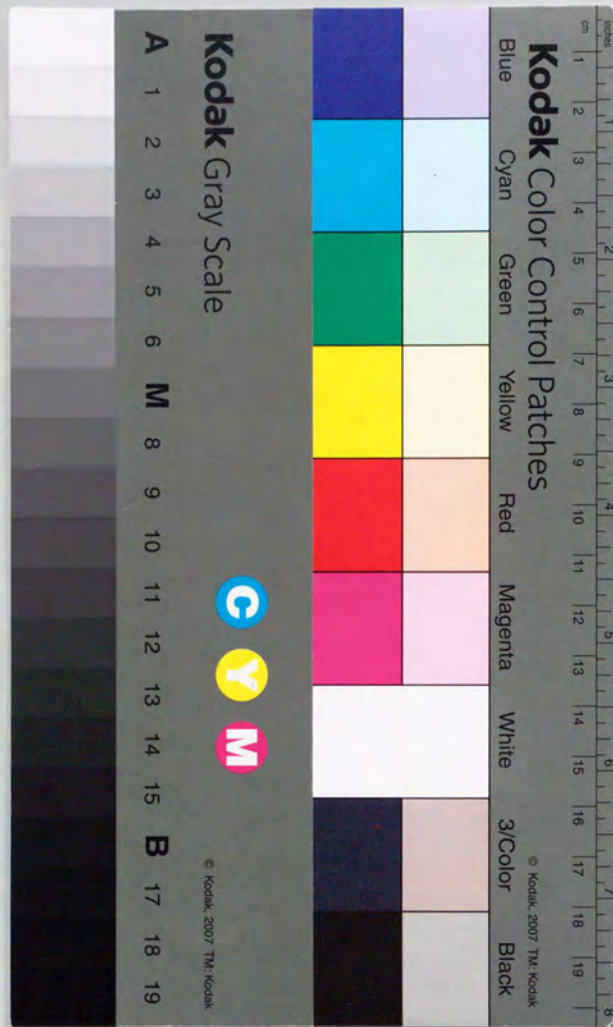


Studies on Control and Dimensioning  
in Asynchronous Transfer Mode Networks

非同期転送モード網における  
制御と設計に関する研究

Hirosaki Saito

斎藤 洋





①

Studies on Control and Dimensioning  
in Asynchronous Transfer Mode Networks

非同期転送モード網における  
制御と設計に関する研究

Hiroshi Saito

斎藤 洋

## Contents

I	Introduction	5
II	Historical Background	17
	1. Introduction	
	2. Circuit switching	
	3. Hybrid switching	
	3.1 Hybrid switching with variable-length frame	
	3.2 Hybrid switching with SAD	
	4. Burst switching	
	5. Packet switching	
	5.1 Reconstruction of continuous speech	
	6. ATM	
	6.1 Modeling of video cell arrival processes	
	6.2 Models of superposed process and general burst traffic	
	7. Conclusions	
	Appendix 1. Speech activity model	
	Appendix 2. Flow control with bit reduction	



### III The Departure Process of an N/G/1 Queue

57

1. Introduction
2. Model and notation
3. Preliminaries
4. The interdeparture times of an N/G/1 queue
5. The interdeparture times of an N/D/1 queue
  - 5.1 Mean interdeparture time
  - 5.2 Second moment of the length of interdeparture times
6. Counting process
7. Examples
  - 7.1 Departure process from an M/M/1 queue
  - 7.2 Departure process from an MMPP/D/1 queue
  - 7.3 The departure from a packetized voice multiplexer
8. Conclusions
- Appendix.  $C_p(z)$  characterization

### IV Optimal Queueing Discipline for Real-Time Traffic at ATM Switching Nodes

82

1. Introduction
2. Model
3. The optimal queueing discipline
4. Implementation
5. Cell loss probability of each class
6. A numerical example
7. Conclusions
- Appendix. Proof of the proposition

### V Optimal Control of Variable Rate Coding in ATM Networks

103

1. Introduction
2. Problem formulation
3. Structure of the optimal control
4. The optimal parameters
5. An average waiting time constraint
6. Conclusions
- Appendix 1.
- Appendix 2.
- Appendix 3.
- Appendix 4.

### VI Optimal Control of Variable Rate Coding with Incomplete Observation in ATM Networks

130

1. Introduction
2. Problem formulation
3. Optimality equations
4. Suboptimal control
5. Performance of suboptimal control
6. Conclusions
- Appendix.

### VII Call Admission Control in an ATM Network without Using Traffic Measurement

164

1. Introduction
2. Preliminaries
3. Upper bound of cell loss probability from MNA and ANA
4. Upper bound of cell loss probability from ANA and VNA



5. Implementation of call admission control	
6. Numerical examples	
7. Conclusions	
Appendix 1.	
Appendix 2.	
Appendix 3.	
VIII A Simplifying Dimensioning Method of ATM Networks	197
1. Introduction	
2. Preliminaries	
3. Dimensioning for a single class of traffic	
4. Dimensioning for multiple traffic classes	
4.1 Algorithm	
5. Numerical examples	
5.1 Example 1	
5.2 Example 2	
6. Conclusions	
Appendix 1.	
Appendix 2.	
Appendix 3.	
Appendix 4.	
IX Conclusions	226
Acknowledgements	230
References	231

## Chapter I

### Introduction

Recently, there has been considerable interest in integrated networks. This is partly the result of the specification of international standards for the Integrated Services Digital Network (ISDN) and its evolution. The ISDN standards define different classes of network services and a single integrated interface for access to these services, although these services may be provided in more than one way.

Advances in communications technology are another motivation for integrating telecommunication networks. Digital transmission is replacing analog transmission, because of decreasing costs, increasing data traffic, increasing capabilities for integrating services, and also to provide higher quality services. The introduction of fiber optics has reduced transmission costs and paved the way for new services employing wide bandwidth.

Furthermore, combining telecommunications services into the same network is widely recognized as offering several benefits. Lower equipment and communications costs can



be achieved through higher utilization as a result of integration, flexibility for introducing and distributing new services can be expected to increase, and terminals will become more portable owing to the integrated interface. Large scale production of highly integrated system components for a unique ISDN will lead to cost-effective solutions.

ISDN is conceived to support many kinds of future services, including broadband services [Händel 89, Tominaga 89] such as those listed below.

- broadband video telephony and video conference
- video surveillance
- high-speed unrestricted digital information transmission
- high-speed file transfer, teleaction and telefax
- video and document retrieval service
- TV distribution (existing, extended, and high definition quality)

The asynchronous transfer mode (ATM) is a target technology in achieving the broadband ISDN (B-ISDN).

The 'circuit switching' and 'packet switching' are widely used transfer modes. However, low utilization of network resources and inflexibility in available bandwidth are disadvantages of the circuit switching transfer mode, and delay-related problems are characteristic in the packet switching transfer mode.

To overcome these problems, many approaches have been proposed: fast circuit switching, enhanced circuit switching, burst switching, hybrid switching, and simplified packet switching protocols. These approaches all suffered from disadvantages, and ATM was proposed. ATM can convey multi-media information, including voice and video [Turner 86b,

Luderer 87, Kawarazaki 88, Händel 89] (Figure 1). Its features are:

- Information is sent in short fixed length blocks called cells. Flexibility to support a variable transmission rate, is accomplished by transmitting the necessary number of cells per unit of time.
- The principle is a basic low-layer function that involves no complex flow control on a link-to-link basis. Flow control and error correction are proposed on an end-to-end basis as needed.
- The switch is self-routing and implemented by hardware. Each cell is individually identified and processed on the basis of a virtual circuit.
- The ATM header contains the label called Virtual Channel Identifier and multiplexing is done by means of labels.

The following sequence is a typical example of communication in an ATM network.

- (i) On requesting connection of a call, the network judges whether network resources along an appropriate route can be (statistically) allocated, based on the anticipated traffic characteristics and the requested grade-of-service (GOS) of the connection request, and the network load. Here, the anticipated traffic characteristics is estimated from the traffic parameters specified by the connection request. The connection request is rejected, when network resources cannot be allocated. This judgement is called call admission control.
- (ii) When the connection request is accepted through call admission control, cells are generated according to the information of the call and are transmitted (Figure 2).
- (iii) Transmitted cells are controlled to maintain cell-level GOS standards, if necessary



(cell level control). The traffic flow of the call is monitored and by policing control, to ensure that it conforms to the specified traffic parameters.

ATM based B-ISDN will offer the following benefits.

- Adaptable to new services with different bandwidth requirements.
- Can integrate circuit switched networks and packet switched networks.
- High utilization of links, because of statistical multiplexing of bursty traffic; also ATM networks, in which multiplexing is done by means of labels, need not consider digital hierarchy.

Although ATM is very promising, traffic design and traffic control present problems.

#### (1) Traffic design

Traffic design includes network architecture design and network dimensioning. The effects of introducing ATM technology for network architecture are, change of costs in nodes and links, dynamic network reconfiguration capability through ATM cross-connects, and improvement of transmission link utilization because of label multiplexing. Thus, the network architecture of ATM networks can be expected to be simpler than synchronous transfer mode (STM) networks [Sato 89].

Dimensioning of ATM networks is another problem. As explained in Chapter VIII, the dimensioning of ATM networks is divided into two levels: the call level and the cell level [Hui 88, Filipiak 89, Saito 89c]. Call-level dimensioning, which provides the number of virtual circuits, can employ the dimensioning method developed for STM networks. The difficulties lie mainly in cell-level dimensioning, which provides output buffer size and output link (path) capacity. In particular, under heterogeneous traffic characteristics

and GOS standards, it is necessary to develop cell-level dimensioning. However, cell-level dimensioning becomes a difficult problem under these conditions.

#### (2) Traffic control

Traffic control can also be divided into the call level and the cell level. (Network reconfiguration through ATM cross-connects in the path level may be added to these level controls [Burgin 89a].)

Call admission control, bandwidth reservation, and routing are controls on the call level. Call admission control with bandwidth reservation must decide whether to accept a new connection according to knowledge of the current network loading, the new connection's anticipated traffic characteristics, and its GOS standards. The determination of an appropriate minimum set of user-specified traffic parameters which yields the anticipated traffic characteristics of the user, and the method of judging whether a new connection can be accepted or not according to traffic parameters are challenging problems.

Routing is another problem in call-level control. However, it is not clear whether ATM networks require special routing schemes.

Cell-level control aims to adjust cell-level GOS performance and to monitor and enforce the cell stream to ensure that it conforms to its specified parameter values. The latter called policing. The policing mechanism depends on traffic parameters, and remains unestablished. The candidates for policing mechanism are the leaky bucket [Turner 88], cell and timer counter [Kowalk 88], and virtual leaky bucket [Gallassi 89]. The remainder of this thesis assumes that the policing control functions well and that traffic actually satisfies specified traffic parameter.



This thesis studies teletraffic issues in ATM networks. Two approaches to control and dimensioning are fundamental to this thesis; a synthetic approach, and a non-parametric approach.

*Synthetic* means that traffic control is synthesized given performance measures and statistics of traffic, from the traffic point of view. Present approaches of traffic studies to these problems are analytic; i. e., traffic control schemes are proposed by an experienced engineer, alternatives are compared by traffic analysis, and dimensioning gives the amount of resources to achieve GOS standards. This new synthetic approach attempts to directly synthesize traffic control which is expected to achieve given performance standards. This approach considers performance measure and statistics of traffic, first. Hardware configuration is considered after control structure is synthesized, although it is considered first in the conventional approaches. This approach may create a quite novel control scheme.

The *non-parametric* approach is explained, by comparing it with the existing traffic studies. Adopting conventional traffic approaches for ATM networks has several disadvantages. One is that higher moments of cell interarrival times and/or correlations among them must be required to be measured, because cell arrival models contain many parameters. (Simple models not using higher moments sometimes overestimate performance and sometimes underestimate it.) Also, results are based on assumptions, such as that the length of silence periods of voice traffic is exponentially distributed, which cannot be supported or validated from field data. Another disadvantage is that while ATM networks will be flexible for introduction of new services, the existing traffic studies require modeling of cell arrival processes for individual services. This reduces flexibility in traffic design and

control for new services. Non-parametric traffic engineering attempts to eliminate these disadvantages. A design method and control should be established that does not require modeling and assumptions and use only items that are easy to measure.

The remainder of this thesis is organized as follows.

#### [Chapter II]

Chapter II reviews traffic studies concerning the integration of networks. This chapter covers circuit switched, hybrid switched, burst switched and packet switched integrated networks, and also ATM networks. The emphasis is on performance considerations for integrated networks.

#### [Chapter III]

It is hoped that, by concatenating several nodes in ATM networks, burstiness may be reduced. To examine this phenomenon, the cell arrival processes at the input and output sides of a node are required to be compared.

In Chapter III, the departure process of an  $N/G/1$  queue is investigated. The  $N$ -process, a versatile point process, can model video/voice cell arrival processes [Yamada 89, Saito 91]. Thus, the smoothing effect of passing through a node can be quantitatively evaluated with the theory developed in this chapter.

In [Saito 91], the results presented in this chapter are directly applied and the following conclusion is obtained: When transmission efficiency is low, which is the case when video traffics GOS standards are fulfilled, burstiness is not reduced by passing through nodes. Consequently, the argument for one node can be applied to all nodes in a network, and it is sufficient to consider a single node.



[Chapter IV]

Chapters IV-VI synthesize cell level controls. Delay quality control in Chapter IV is a queueing discipline which deals with delay-sensitive traffic. The optimal discipline is derived which minimizes the number of cells delayed beyond the maximum allowable time specified for an individual call class, without assumptions concerning the cell arrival process and buffer management schemes. Implementation of the optimal discipline is discussed.

[Chapter V-VI]

In Chapters V and VI, the optimal control for selective cell discarding is derived. Under embedded coding of voice signals, cells of sampled voice signals have different significance. That is, one cell may contain bits of higher significance than another cell. When a network is congested, the bits are dropped by discarding cells at vocoders or multiplexers located at the network entry point and at the ATM switching nodes within the network. Thus, there is a tradeoff between deterioration in voice quality and reducing congestion. The optimal control is shown to maximize the long-run average coding rate under an average cell queue length constraint.

In Chapter V, the numbers of cells arriving at different slots are mutually independent. This is the case when sufficiently many sources are multiplexed. When the number of voice sources is small, the numbers of cells arriving at slots are not independent. In this case, some mechanism forecasts the number of arriving cells by means of estimating the number of active voice sources (Chapter VI). Here, 'active' denotes the off-hook user in talkspurts.

The synthesized optimal control in Chapters V and VI is shown to have simple structure, and to be the random selection between two feedbacks. In particular, when voice cells

are divided into two levels of significance, that is, most significant and least significant, the optimal control is bang-bang control such that when a queue length exceeds a threshold, the least significant cells are discarded.

In Chapters V and VI, voice and data traffic are assumed. However, this can be extended to the case such that embedded coded video traffic is added to voice and data traffic.

[Chapter VII]

In Chapter VII, a call admission control based only on parameters specified by users is proposed. The key technology here is non-parametric evaluation of the upper bound of cell loss probability. Non-parametric evaluation is done using only the specified parameters, and other assumptions, parameters and modeling are not needed.

The proposed call admission control rejects connection requests when the evaluated upper bound of cell loss probability exceeds a cell loss probability standard. Thus, the cell loss probability standard is guaranteed to be satisfied under this control. Implementation of this control to quickly evaluate cell loss probability after acceptance of a new call is discussed.

When there is no information on cell arrival processes except for the specified parameters, the result of this chapter is available for use in the dimensioning of ATM networks.

[Chapter VIII]

In Chapter VIII, a 'non-parametric' dimensioning method for ATM networks is described. The dimensioning method here is divided into the call level and the cell level. Call-level dimensioning provides the number of virtual circuits, considering the call loss



probability standard. Cell-level dimensioning yields output buffer size and output link capacity, considering the cell loss probability standard and the maximum admissible delay. This dimensioning method is applicable to multiple traffic classes. It employs the probability density function of the number of cells arriving in a fixed interval. However, the p.d.f. need not be parameterized. The measured frequency distribution can be used as the p.d.f. This is not the case in the ordinary dimensioning method, which is based on queueing theory.

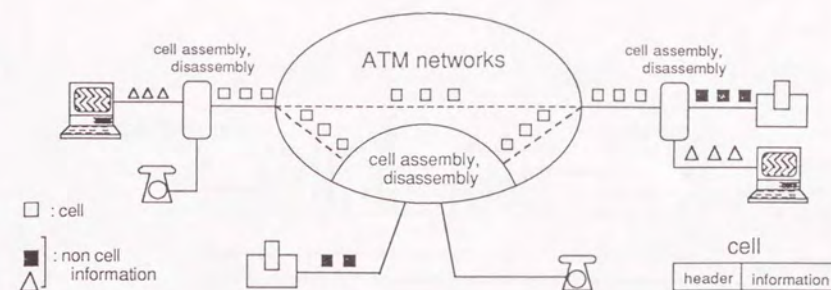
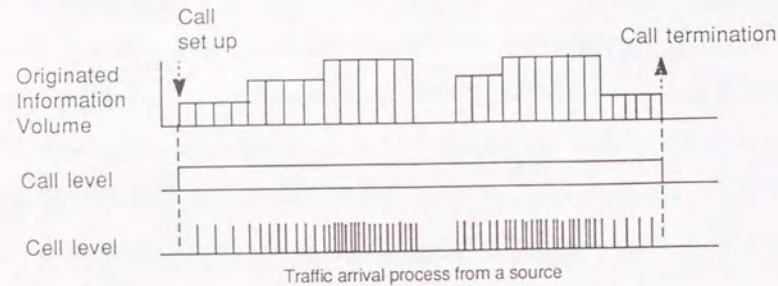


Figure 1. ATM networks

- User data (voice/data/video) is transferred through the networks in short fixed length blocks called cell.
- Networks provide the basic functions to transfer cells, which are required by each service.
- Functions dependent on services, e.g. cell assembly/disassembly, flow control and error recovery, are supported by the interface between networks and their outsides or terminals.





	Traffic descriptor	GOS definition	Transfer mode
call level	<ul style="list-style-type: none"> <li>• Traffic intensity (call/sec)</li> <li>• Offered load (erl.)</li> </ul>	<ul style="list-style-type: none"> <li>• Call set-up delay</li> <li>• Call loss prob.</li> </ul>	<ul style="list-style-type: none"> <li>• Virtual circuit set-up</li> </ul>
cell level	<ul style="list-style-type: none"> <li>• Cell intensity (cell/sec)</li> <li>• mean, variance?</li> </ul>	<ul style="list-style-type: none"> <li>• Cell transfer delay</li> <li>• Cell loss prob.</li> </ul>	<ul style="list-style-type: none"> <li>• Cell transfer</li> </ul>

Figure 2. Parameters for ATM network

## Chapter II

### Historical Background

*Interest in integrated communications networks has been stimulated recently. This is partly the result of the specification of international standards for the Integrated Services Digital Network (ISDN) and its evolution. Advances in communications technology are another motivation for integration of networks. Furthermore, combining telecommunications services into the same network is widely recognized to offer several benefits. Hence, many efforts have been made and many approaches have been proposed for integrating telecommunication services.*

*This chapter addresses the issues of transmitting voice and data in public integrated networks. The emphasis is on reviewing performance considerations for integrated networks. Teletraffic issues for video traffic are also discussed. [Saito 90c]*



## 1. Introduction

Today's public telecommunications networks include at least two types of stand-alone network. One network is dedicated to voice and the other network is dedicated to data communications. Each network has its own interfaces, transport system, and switching mechanism, which are not suitable for other types of applications.

The characteristics of voice and data signals are fundamentally different, which is also the reason for dedicated networks. Voice signals are inherently real-time analog signals generated by human speakers. Voice traffic can tolerate a certain percentage of errors without becoming objectionable. However, a requirement for a delay in networks are severe and the maximum allowable delay is in the range of 40 to 250 ms [Gruber 83, Turner 83, 86a, Green 87]. On the other hand, most data is machine-generated and digital. A subscriber with interactive data traffic can wait a sizable fraction of a second, and bulk data traffic between machines can wait much longer. Data traffic cannot tolerate errors. Data traffic communication is generally asymmetric and may be quite bursty [Chen 88].

Interest in integration of communications has been stimulated recently. This is partly the result of the specification of international standards for the Integrated Services Digital Network (ISDN) and its evolution. The ISDN standards define the different classes of service of a network and a unique integrated interface for access to these services, although these services may be provided in more than one way.

Advances in communications technology are another motivation for integrating voice and data networks. The telephone system is gradually being converted from an analog to an entirely digital network, because of decreasing costs, increasing data transmission,

increased capability for integrating services, and in order to provide higher quality services. The introduction of fiber optics has reduced transmission costs and permitted new services employing wide bandwidths.

Furthermore, combining telecommunications services into the same network is widely recognized to offer several benefits. Lower equipment and communications costs can be achieved through higher utilization as a result of integration, flexibility for introducing and distributing new services can be expected to increase and terminals will become more portable owing to the integrated interface.

Hence, many efforts have been made and many approaches have been proposed for integrating telecommunication services. Each technique has its own characteristic features, advantages and disadvantages.

This chapter addresses the issues of transmitting voice and data in public integrated networks by circuit switching, hybrid switching, burst switching and packet switching. The emphasis is on reviewing performance considerations for integrated networks. In the section of ATM, teletraffic issues for video traffic are also discussed.



## 2. Circuit Switching

In circuit switching, a complete end-to-end circuit is established for each pair of users and dedicated for the full duration of use. The characteristics of circuit switching are: access blocking, call setup/clear time, fixed throughput, and low delay after setup [Gruber 81b]. These characteristics are considered acceptable for voice. However, for interactive data applications, low utilization and call setup have been considered as disadvantages of circuit switching.

Fast circuit switching and enhanced circuit switching overcome these disadvantages [Harrington 80, Gruber 81b]. In fast circuit switching systems, the signaling speed to set up and break a connection is so fast that circuits are not dedicated to the interactive data user during thinking time.

Enhanced circuit switching is accomplished by using traditional circuit switching supplemented with TASI (time assigned speech interpolation) [Bullington 59] or DSI (digital speech interpolation) [Campanella 76] for voice traffic and adaptive data multiplexing which manages data on a loss or blocking basis. Enhanced circuit switching as well as fast circuit switching can improve the utilization of a transmission link by dedicating it to a user only while he is sending information. Both techniques require a device which can allocate a circuit to a user very rapidly.

Akiyama et al. evaluated an integrated system similar to enhanced circuit switching [Akiyama 86]. They analyzed the state equations for the number of voice calls and the number of data calls in the system, assuming that voice and data calls arrive in Poisson processes and that their holding times are exponentially distributed. Both voice calls and

data calls are managed on a loss basis. A voice call requires  $m$  times as much bandwidth as a data call. Numerical examples indicate that the utilization of circuits increases by 8-18% because of integration.

In spite of many efforts, however, circuit switching technologies have several deficiencies associated with providing integrated voice/data. For example, there is no error detection or correction. Thus the future lies in effective use of packet and hybrid switching technologies [Harrington 80]. Hence, interest has recently been focussed on fast packet switching and so on.

## 3. Hybrid Switching

A natural concept for combining voice and data communication is a hybrid switch, which has been actively studied since the 1970s. A hybrid switch uses a TDM (time division multiplexing) frame structure. In TDM, frames are subdivided into slots, and are partitioned by boundaries into two categories: circuit-switched traffic, including voice, and packet-switched traffic, including interactive data. Circuit-switched traffic is managed on a loss basis and packet-switched traffic is managed on a delay basis. A boundary can be fixed or movable, and a movable boundary enables packet data to seize currently idle circuit voice slots during a particular frame. Note that a variety of enhancements pertaining to the basic TDM frame are possible, including fixed- or variable-length frames, fixed or variable numbers and sizes of slots, and fixed or variable numbers and positions of boundaries [Bially 80b, Gruber 81a, b, Ross 82, Green 87] (Figure 1).

A typical example of a hybrid switch is a SENET [Coviello 75]. Performance has been determined by exact analysis [Fischer 76, Kwong 84]; approximate analysis [Kümmerle 74,



Lehoczky 81a, Gaver 82, Leon-Garcia 82, Honig 84]; and simulation [Weinstein 80]. In addition, [Yum 87] investigates the effectiveness of routing, and [Maglaris 82, Avellaneda 82] and [Shioyama 85] optimize the allocated number of voice slots.

Although the fluid flow approximations [Gaver 82, Leon-Garcia 82, Honig 84] and the diffusion approximations are computationally attractive, modeling and exact analysis in [Kwong 84] is shown here. Voice (circuit-switched traffic) and data (packet-switched traffic) arrive in two independent Poisson processes. Calls arriving during a frame must wait until the beginning of the next frame. Voice traffic is allowed to queue for the duration of one frame but no longer. If the number of free voice slots is greater than the number of voice calls ahead of it in the buffer, the call receives service; otherwise, it is lost and leaves the system. The holding time for voice calls has a negative exponential distribution. For the data traffic, arrivals are buffered and, at the beginning of every frame, are placed in the data and unoccupied voice slots (movable boundary hybrid switching). The data traffic remains in the system and is not lost if it does not receive service. Data packets are assumed to be of a constant size equal to the duration of a frame. Buffers are assumed to have an infinite capacity. Voice and data slots are modelled as servers, and the queueing model is analyzed, using the two-dimensional imbedded Markov chain given by the joint process on the numbers of voice and data calls in the channel immediately after the beginning of the frame. The generating functions are determined by the transition probabilities of the imbedded Markov chain and the average data delay is derived. The determination of the average data delay requires finding the roots of polynomials.

### 3.1 Hybrid switching with variable-length frame

A typical SENET concept has fixed frame- and slot-lengths. In a hybrid multiplex structure with a constant frame format, such as a SENET, either packets are too long for transmission in that frame, in which case capacity is wasted, or the packet length is shortened, which increases the overhead of that packet relative to its information field. If both frame and packet are allowed to have variable length, then unused time intervals in the communication link can be eliminated. Miyahara and Hasegawa [Miyahara 78] proposed hybrid switching with a variable length frame and packet. They evaluated its performance and championed its effectiveness. In [Maglaris 82], and [Janakiraman 84a] as well as [Miyahara 78], performance of a hybrid switch with variable frame length is discussed. The average data delay and the average transmission time in [Maglaris 82], and the channel utilization, the voice blocking probability and the average data delay in [Janakiraman 84a] are derived.



### 3.2 Hybrid switching with SAD

In a SENET system, a data queue builds up significantly and its size becomes unacceptably large [Bially 80b, Weinstein 80, Green 87, Li 88a]. This is because, even though the system is capable of handling all of the offered traffic on average, there are long periods of time during which the voice traffic occupies most of channels capacity.

One approach to overcome this disadvantage is to reserve data capacity [Okada 86, Lehoczky 81a, Gaver 82]. However, since the holding time of voice traffic is extremely long compared with that of data traffic, data traffic can be transmitted only via reserved capacity during long periods; thus a large reserved capacity is necessary for an acceptable data delay. Therefore, flow control should be imposed in a SENET.

A promising technique for removing this problem is speech activity detection (SAD, see Appendix 1). A hybrid switch with SAD assigns a slot only for an activated voice call. The silence in a voice is interpolated to be either data or a talkspurt of another voice. If a talkspurt cannot be allocated to an idle slot, it is frozen out, resulting in clipping. The relationship between speech clipping and subjective speech quality is shown in [Gruber 85].

Since talkspurt/silence variations are more rapid than those due to call initiation and termination, data queues build up and discharge over shorter time intervals and average queue lengths are shorter. A hybrid switching with SAD is discussed in [Fischer 79, Bially 80b, Gruber 81b, Lehoczky 81b, Gaver 82, Ross 82, Arthurs 83, Sriram 83, Williams 84, Li 85, O'Reilly 85,86b, Ahmadi 86, Green 87, Vakil 87], and [McDysan 88].

The first attempt to analyze this system was by [Fischer 79]. A single channel system

is analyzed by the steady state equations and the moment generating functions about the state of a voice call and the number of data packets on the assumption that the silence period has a hyperexponential distribution. Performance measures including the mean waiting time of data packets are obtained, assuming that talkspurt length is exponentially distributed.

Techniques for analysis used in other works are the matrix geometric method [Neuts 81] in [Lehoczky 81b], and [McDysan 88]; numerical techniques for the steady state equations in [Williams 84]; discrete time queuing analysis in [Sriram 83], and [Li 85]; and a fluid flow approximation in [O'Reilly 85,86b], and [McDysan 88].

A queueing model of hybrid switching with SAD is very similar to that of a burst switch. The difference is that in a burst-switched system the frame capacity is channeled and the transmission speed of data traffic is the same as that of voice traffic [Lim 86]. Hence, if the voice transmission speed is assumed to be equal to data transmission speed and the data packet length is assumed to be equal to the frame duration, then either analysis will also give a rough approximation of the results of the other analysis.

Let us consider the buffered speech interpolation in a hybrid switching as an advanced version of a hybrid switching with SAD [Weinstein 79, Gruber 81b]. In buffered speech interpolation systems, talkspurt delay occurs instead of freezeout, and utilization of a transmission link can be higher than in fractional loss systems, particularly when a small number of users is multiplexed [Weinstein 79, Janakiraman 84b].

Analytical results for this system are derived in [Konheim 84], [Liang 85], [Kim 88] and [Luhanga 88]. (In [Fischer 80] buffered speech interpolation without data integration



is considered. A numerical technique based on the state equations gives an approximated analysis.) In [Konheim 84], the performance of a voice/data multiplexer using a movable boundary is considered under various allocation schemes of slots for voice and data. The analysis is based on the joint process  $(v, d)$ , where  $v$  is the number of voice packets in the system and  $d$  is the number of data packets in the system. The number of allocated slots for voice packets is determined by allocation function  $a(v, d)$ . A numerical technique that is a simplified version of [Konheim 84] is given in [Liang 85]. Allocation schemes based on a priority discipline are considered in [Kim 88], which is available as a study of queueing disciplines in a slotted packetized multiplexer. A fluid flow approximation following Anick's model [Anick 82] is employed in [Luhanga 88]. Since it is supposed that voice traffic has preemptive priority over data, the analysis can be simplified by considering the voice process first and then analysing the data process. Studies by [Li 88a,b,c] provide useful additional information.

Additionally, Hayashida et al. evaluated the flow control method whereby data traffic is regulated according to the number of voice calls in the system [Hayashida 86].

#### 4. Burst Switching

Burst switching, a method for switching voice and data in an integrated way, has been proposed since the 1980s [Amstutz 83, Haselton 83]. It is a form of message switching that combines different features from the fast circuit switching. A burst begins with a header, which contains the network address of the burst's destination and information describing the burst type (voice/data/command). The burst header is followed by information with a completely variable length, that is, a talkspurt or data message. A burst ends with a burst termination character called FLAG. Burst switching uses headers for routing and bandwidth contention queueing (Figure 2).

A burst is sent between switches through a TDM channel, and at the completion of transmission, the channel becomes available for reassignment to another burst. Characters of a burst arriving at a switch are buffered and the buffer is placed in the an appropriate link output queue as soon as sufficient information for routing is available. In other words, a burst can begin to be forwarded before it is buffered completely. That is, the burst switching is not of the store-and-forward type.

Channel congestion occurs if there are more bursts in a link queue than there are idle channels in the link. There are three priorities; high, normal and low for assignment of a channel. Command bursts have high priority, voice bursts have normal priority, and data bursts have low priority (Figure 3). If 2ms of voice samples have been accumulated and output has not begun, the accumulated characters of the voice burst are discarded. This causes clipping.

Simulation studies for burst switching performance are reported in [Morse 85], and



[Jack 86]. Network topologies are mainly discussed in [Morse 85], and the effectiveness of a high speed "superchannel" (virtual channel comprising several channels temporarily linked) for bulk data traffic is discussed in [Jack 86]. In these results, the voice traffic model is based on the empirical results in [Yatsuzuka 82a] (see Appendix 1).

Exact analyses are given in [Descoux 85, Lim 86], and [Aboul-Magd 88]. In [Descoux 85], bursts are assumed to arrive in a Poisson process and the effect of buffering of voice bursts is evaluated. In [Lim 86] and [Aboul-Magd 88], the lengths of talkspurt and silence periods are assumed to be exponentially distributed, the arrival stream of data bursts is assumed to be Poisson, and the duration of data messages is assumed to be exponentially distributed. The number of off-hook voice sources is fixed. An infinite buffer is considered in [Lim 86] and a finite buffer is considered in [Aboul-Magd 88]. The admissible waiting time of voice bursts is neglected in their analysis, that is, voice traffic which cannot find an idle channel begins to be discarded immediately until one can be found. (This assumption is commonly used in many performance studies on burst switching.) The joint process (the number of data bursts in the system, and the number of voice calls in an active mode) is considered and solved by a matrix geometric method [Neuts 81] in [Lim 86] and by a computational procedure for the state equations in [Aboul-Magd 88]. The average queueing delay of data bursts is numerically evaluated.

Approximate analyses using the same assumptions on the lengths of talkspurt/silence, the arrival process, and the holding times of data bursts are reported in [Ma 87,88], and [Aboul-Magd 88]. In [Ma 87], a delay scheme is introduced in which a voice burst that cannot enter the service upon arrival is allowed to wait in a queue for up to a specified

time before part of the voice burst is discarded.

In [Aboul-Magd 88], a quasi-static approximation is derived, in which the process of the number of data bursts in the system is assumed to be stationary for each state of voice bursts. A fluid flow approximation which assumes that the silence periods of a voice source are hyperexponentially distributed is given in [O'Reilly 87a]. The results show that the hyperexponential assumption, which is supported by empirical data [Yatsuzuka 82a], is crucial in data performance, and the assumption that the silence duration is exponential provides an overly optimistic result.

An analytical result for a burst switched network is given in [Leon-Garcia 86]. The author derives a fluid flow approximation for a tandem network without making an independent assumption. Infinite data buffers are assumed. A data burst is served at the second stage of the tandem queues after the service completion at the first stage. Voice bursts occupy servers at the first and second stages, or occupy a server at the second stage. Linear first order partial differential equations for the number of data bursts in each stage are derived.

Studies between burst and fast packet switching are compared in [Li 87] and [O'Reilly 86a,87b]. Transmission in burst switching is characterized by talkspurt clipping and the clipping can be large although only a few talkspurts are clipped. In packet switching, it is characterized by a packet delay and the delay for each packet is small, although many packets may be delayed. Therefore, better performance can be achieved in packet switching, given the assumption that a 20-ms worst-case mean clipping period in burst switching is equivalent to a 20-ms worst case mean packet delay in packet switching [Li 87].



If we take 2% packet loss to be equivalent to 0.5% freezeout, and follow O'Reilly's analysis, using his earlier results on burst switching performance [O'Reilly 85,86b] and Jenq's result [Jenq 84] for packet switching performance with Stanford's approximation [Stanford 85], packet switching can provide the same performance as burst switching with the same TASI advantage and equal capacity. Although the average residual capacity left for data is significantly greater for burst switching than for packet switching, data performance within that residual capacity is not significantly different [O'Reilly 87b]. Combining these results, it can be concluded that both techniques, burst and packet switching, have roughly equivalent performances for voice and data.

Table 1 summarizes the mentioned-above result.

## 5. Packet Switching

The integration of digital voice with data in a packet switched network has been tried many times since the initial experiments on ARPANET in 1974. A series of packet speech systems experiments has been conducted under the sponsorship of the Defense Advanced Research Projects Agency [Weinstein 83]. Packet switching provides a powerful mechanism for dynamically sharing transmission resources among users with time-varying demand and can be expected to achieve high utilization of transmission links. Hence, packet switched networks were originally suitable for applications requiring low throughput and low delay (e.g. interactive data), and those requiring high throughput and accepting higher delay (e.g. file transfer). However, packetized voice requires both high throughput and low delay, which is not consistent with the capabilities of packet-switching techniques. Delay related issues are particularly characteristic in packetized voice system. See [Gruber 85] for the

relationship between delay and subjective speech quality.

To overcome some of the delay related problems associated with packetized voice, various advanced design concepts for packet switching have been considered. Shorter packets, which are used to reduce packet assembly delay, result in excessive overhead and subsequent inefficient resource utilization. Therefore, abbreviated headers are required, and these are made possible by virtual call routing. Packet voice protocols that have time stamps, speech playout with buffering and no error control have been proposed [Gruber 81b, Weinstein 83, Coviello 79].

In addition, there are movements to achieve integrated services networks including data, voice and video communications, by a packet technique called fast packet [Turner 83,85,86a,b, Kirton 87, Giocelli 87, Luderer 87, Bris 86, Huang 84]. The key technologies that make the integrated services packet networks possible are high speed digital transmission facilities with excellent error performance, simple link level protocols, and hardware implementation of basic switching.

According to [Turner 86a], current plans for ISDN relying on circuit switching, which requires that the available bandwidth be divided up into fixed-size channels, are too inflexible to satisfy various requirements in the future. Furthermore, integrated services packet networks have advantages. Packet switching systems can adapt to changing demands, such as new services and new technologies, can provide both an integrated customer interface and a single network solution for a wide range of communications needs, leading to substantial cost savings in switching systems and system administration, and can exploit communications burstiness which makes poor use of conventional circuit switched facilities.



Since packet switching has recently been recognized as a promising technique, there have been a considerable number of studies that analyze the characteristics or measure performance in packet voice/data networks (See Table 2). These analyses all assume that packets are generated at regular intervals during a talkspurt and that no packets are generated during a silence period (Figure 4). They also assume that the lengths of talkspurts and silence periods in a single voice source are exponentially distributed, except for a simulation analysis in [Seguel 82], which uses an experimentally derived voice model (See Appendix 1). The length of voice packets is fixed when there is no flow control (See Appendix 2). The remaining results in this section are based on the use of speech activity detection (SAD). However, [Gopal 86] indicated that not using SAD has advantages: a simpler playout strategy can be used, and 32 Kbps without SAD provides comparable or better performance than 64 Kbps with SAD. Reduced transmission costs resulting from the introduction of optical fibers, the use of embedded coding (see Appendix 2) and the requirement of high quality voice communication tend to reduce the advantages of SAD. Therefore, not using SAD will be a promising alternative to using SAD in the future.

Optimal packet lengths are analyzed in [Minoli 79, Suda 84], and [Giorcelli 87]. The packet length which minimizes the total delay (= queueing delay - packetized delay) is 300 to 700 bits, assuming that packets arrive in a geometric distribution, the speed of transmission link is 50 Kbits/s, and the header is 50 bits long [Minoli 79].

A simulation study [Suda 84] shows that the optimal packet length is 250 to 300 bits when the speed of transmission link is 1.544 Mbits/s. An experiment using a Delta network shows that a packet length of 12 bytes is appropriate [Giorcelli 87]. These results indicate

that short packets are preferable in integrated voice and data packet networks, whereas in fact long packets are used in conventional networks.

Performance evaluation for a voice/data packet multiplexer or a output queue of a transmission link in integrated voice/data packet networks which receives much attention, has difficulty in modeling the superposed voice arrival processes (Figure 5).

The aggregated arrival process is highly correlated, so the accuracy of the simplest approximation, the M/D/1 model, for the voice/data packet multiplexer is either too poor for practical use ([Seguel 82, Daigle 86]) or is acceptable only if the number of voice sources is large [Kim 83], traffic is light [Jenq 84] or the number of waiting rooms is small [Sriram 86].

Therefore, the analysis is made tractable by using a renewal process approximation [Jenq 84, Sriram 86]; a Markov modulated Poisson Process (MMPP) approximation [Heffes 86, Daigle 86, Ide 88]; a fluid flow approximation [Daigle 86, Li 88d, Tucker 88]; a discrete time approximation [Tanaka 82, Jenq 84, Li 88d] or a semi-Markov process approximation [Daigle 85,86, Stern 83,84].

[Stern 83,84] considers the overload, that is, the states in which the aggregated input stream from active talkers exceeds the transmission capacity of the output link, and the underload, that is, the states in which the aggregated input stream from active talkers does not exceed the transmission capacity of the output link. Using the fact that the point at which an overload/underload cycle begins is an imbedded Markov point, Stern derived a functional equation for the steady-state probability generating function of the number of packets in the system. This equation can be solved by functional iteration and spectral



factorization. Unfortunately, Stern's results are not for stochastic equilibrium, but are for the beginning of an overload/underload cycle.

[Daigle 85,86] analyzed a similar model to [Stern 83,84] by matrix geometric techniques [Neuts 81], which are computationally more attractive than the techniques employed in [Stern 83,84]. Daigle et al. [Daigle 86] also tried to analyze the voice multiplexer by an MMPP/M/1 queueing model and by a fluid flow approximation which yields linear ordinary differential equations on the buffer content and the number of active talkers. This approximation was used in [Anick 82, Gaver 82, O'Reilly 85], and [Tucker 88]. Numerical results in [Daigle 85] and those obtained by fluid flow approximation are both more accurate than the results given by the MMPP/M/1 model. A fluid flow approximation in [Tucker 88] extends the analysis in [Daigle 86] from a fluid flow approximation for an infinite waiting room to that for a finite waiting room.

A matrix geometric approach is used in [Heffes 86] as well as in [Daigle 85]. However, in [Heffes 86], the aggregated arrival process is expressed by the two-state MMPP and, as a result, the calculation of the rate matrix can be derived explicitly. The four parameters of the MMPP are chosen so that the following characteristics of the superposition are matched:

- 1) the mean arrival rate;
- 2) the variance-to-mean ratio of the number of arrivals;
- 3) the long term variance-to-mean ratio of the number of arrivals; and
- 4) the third moment of the number of arrivals.

The characteristics of a single voice source provide the evaluation of these quantities

for the superposition of packet voice processes. Finite buffers and a class of overload control mechanism are also discussed in [Heffes 86].

[Ide 88] also used the MMPP as a model of aggregated arrival processes. Precisely speaking, he used the Interrupted Poisson Process (IPP) as the model of a single source, and employed  $N$  (the number of off hook users) superposed IPPs as an aggregated arrival process and the supplementary variable method to solve the MMPP/G/1. He suggested that the four parameters of the IPP should be matched with mean, variance and peakedness. [Jenq 84] insisted on the verification of the renewal approximation of the superposition of arrival streams in heavy loads, and the queueing model with the renewal approximation arrival process was solved by the queueing network analyzer [Whitt 83].

In [Tanaka 82], numerical evaluation of the steady-state probability using a discrete time approximation agrees well with simulation for mild loads. However, [Jenq 84], who employs a very similar discrete time approximation, shows that it is less accurate than the renewal approximation.

Of the works mentioned above, [Tanaka 82, Stern 83,84, Jenq 84, Daigle 85,86, Ide 88] and [Tucker 88] analyzed a packetized voice multiplexer. [Jenq 84] was extended by [Sriram 86] for an integrated voice and data multiplexer under the assumption of the FIFO service discipline and the Poisson arrival data stream. The other works also can be extended to the case under the same assumption above. Extension is often straight forward. However, another effort is needed to examine the performance of individual classes of traffic (voice and data) for any queueing discipline.

The loss probability of voice packets are analyzed in [Li 88d]. It is assumed that



$N$  voice calls are multiplexed, and that talkspurt and silence periods are exponentially distributed. A two dimensional Markov chain, (the number of voice calls in talkspurt, the number of voice packets in the system), is considered, and attention is paid for the sojourn time in the blocking period. Numerical results are given for a small system by a discrete time two dimensional Markov chain with time epochs set at the end of each frame, and for a large system, a fluid flow approximation is used to obtain numerical results.

A comparative discussion of burst switching and packet switching given in [Li 87] and [O'Reilly 87b] indicates that both techniques performed roughly equivalently for both voice and data. (See Section 4.)

### 5.1 Reconstruction of continuous speech

An important aspect of a packet-switched voice call is the reconstruction of a continuous stream of speech from the set of packets that arrive at the destination terminal, each of which may encounter a different amount of delay in the packet network. Packets are produced at the packet voice sender (PVS) at regular time intervals during a talkspurt and sent through the network. As the packets arrive at the packet voice receiver (PVR), they are reconstructed into a continuous stream of voice samples and delivered to the destination customer.

This reconstruction is done by choosing a target playout time for each incoming packet after the packet is produced. Each packet that arrives before its target playout time is buffered until then to compensate for random network delay. If a packet does not arrive before its playout time, the buffer will be empty, and the late packet is lost or played out immediately, which degrades the speech quality. Therefore, the design of a target playout

time and a playout strategy is an important problem (Figure 6).

[Barberis 80] and [Montgomery 83] state techniques which estimate in some way the delay encountered by each packet to determine how speech is reconstructed. In [Barberis 80], three strategies are stated for the first packet of a talkspurt:

- 1) Null timing information (NTI) device: the PVR delays every first packet of talkspurts by a given amount  $T$ ;
- 2) Incomplete timing information (ITI) device: if the estimated network transit time is below a threshold control parameter  $T$ , then that packet is additionally delayed by an amount equal to the threshold minus the estimated transit time;
- 3) Complete timing information (CTI) device: this device works like the ITI device, but the exact transit time is assumed to be known.

These playout strategies depend on the estimation method of the delay encountered by a packet in the network. No estimation is necessary for NTI mechanism, but the length of a silence period cannot be reproduced exactly because of the stochastic network delay. The ITI algorithm, for example, requires a time stamp which may include a synchronization offset. In the presence of an offset, an optimal voice reconstruction is not possible. The time stamp recorded by the synchronized clock maintained between the PVS and PVR is necessary for the CTI device. Examples of time stamps are given in [Listanti 83, Ueda 83, Turner 83, 86a, Weinstein 83, Muise 86, Green 87, Luderer 87] and [Boyer 87].

Details of synchronization and delay estimation, which is a key technique for the playout strategy, are stated in [Montgomery 83] and [Barberis 80]:

- 1) Blind delay: the PVR makes a worst case assumption about the delay encountered



by a packet;

- 2) Roundtrip estimation: the round trip delay between the PVS and PVR, which is measured by sending a packet between the PVS and PVR, is used to estimate the one-way delay of a particular packet;
- 3) Absolute timing: synchronized clocks are maintained at the PVS and PVR, and packets carry absolute time stamps;
- 4) Accumulated variable delay: the packet network keeps track of the delay experienced by a packet as it travels through the network.

Absolute timing synchronization corresponds exactly to the CTI method, and blind delay to the NTI.

Thus, the performance of each playout strategy has been studied. Suda et al. evaluated mean total delays in the NTI, CTI and NTI-CTI mixed strategies by simulation for the one-hop model and network model [Suda 84]. In [Gopal 84], the performance of the NTI mechanism is evaluated under the 'discard' scheme (a late packet for the target playout time is discarded), and the 'delay' scheme (the playout of a late packet for the target playout time is delayed). The key assumption is that the interarrival times between voice packets in a talkspurt at the destination buffer are independent and identically distributed. Buffers at the destination have infinite capacity. In [Mehmet 86], a design issue of the destination buffer with the output strategies shown in [Barberis 80] for the 'delay' scheme is discussed in a similar way to the delay analysis in [Barberis 80,81].

## 6. ATM

ATM (Asynchronous Transfer Mode) is the specific packet-oriented transfer mode which provides the capability of transmitting multi-media information. It is based on the development of fast packet technologies [Thomas 84, Gonet 86, Boyer 87, Eklundh 88a,b, Kawarazaki 88]. ATM is expected to provide an integrated network which can support voice, data and video, and is the most promising approach to the above-mentioned move toward the integration of networks by packet-switching-oriented technologies. International efforts are being made towards standardizing ATM. Its features are:

- Information is sent in short fixed length blocks called cells. Flexibility, which supports any transmission rate, is accomplished by transmitting the necessary number of blocks.
- The principle is a basic, low-layer function that involves no complex flow control on a link-to-link basis. Flow control and error correction are provided on an end-to-end basis as needed.
- The switch is self-routing and implemented by hardware. Each cell is individually identified and processed on the basis of a virtual circuit.

The simplified protocols and routing result in a concise header format, which enables packets to be processed entirely by hardware.

Fast packet switching for ATM has been realized in a variety of ways: self-routing multi-stage interconnection networks (e.g. STARLITE [Huang 84] and Wideband Packet [Luderer 87]); ring or bus structures with appropriate medium access protocols (e.g. ORWELL ring [Mitrani 86]); and the "Asynchronous Time Division" switching matrix of the PRELUDE System [Gonet 86, Bris 86, Devault 88]. The ORWELL ring is analyzed in



[Mitrani 86] and PRELUDE in [Boyer 87], which models a multiplexer as a  $\sum_i D_i/D/1$  queueing system and a switch as a queueing system with geometric arrivals, considering a pseudosynchronous strategy. The effect of output buffer sharing in the Starlite fabric is evaluated in [Eckberg 88]. Eckberg et al. assume that cells arrive at each link according to a Bernoulli process at each time slot.

Roughly speaking, the queueing model of an ATM switch/multiplexer or output link is similar to that of a packet switch/multiplexer or output link. Therefore, performance analysis for an integrated multiplexer or a output link is available as a rough approximation.

In ATM networks, congestion should be evaluated at a call level as well as a cell level [Hui 88]. Calls should satisfy QOS of a call level, and cells of an accepted call require to satisfy QOS of a cell level. Thus, call admission control and a design method become complicated.

### 6.1 Modeling of video cell arrival processes

In an ATM networks, moving pictures such as video telephony, TV conference, and on-demand video distribution are expected to be a main service. However, statistics on video sources have been studied only since the 1970s [Haskell 72]. In particular, the statistics of a video source of high bit rate coding are a current issue. Bit rate statistics for frame and video source models are summarized in Table 3.

Most statistics are obtained frame-by-frame. Thus, a cell arrival process cannot be directly identified from them. Measurements of the cell arrival process itself are given in [Kishimoto 89] and [Nomura 89a]. [Kishimoto 89] concludes that the cell interarrival time distribution cannot be modeled by a simple distribution, such as an Erlang distributions,

but [Nomura 89a] showed that it can be modeled by the 6th Erlang distribution for a TV program. In general, these results strongly depend on the coding mechanism, and an appropriate model for video traffic has not been established.

The other proposed models of a video cell arrival process are based on frame measurements, and can be divided into those in which the average cell interarrival time within a frame is fixed, and those in which the average cell interarrival time during a certain interval, beginning from the starting point of a frame, is the peak-bit-rate of the coder, and no cells arrive in the remainder of the frame (Figure 7).

In both types, the final mathematical formulation is an MMPP or an auto-regressive (AR) model.

Examples of the fixed average interarrival model are described in [Maglaris 88, Ogino 88, Nomura 89a, Sen 89, Yamada 89]. A continuous Markov chain is employed to model video sources, and statistical multiplexing is analyzed by a fluid flow approximation [Maglaris 88]. An extension which includes scene changes is found in [Sen 89]. In [Yamada 89], an MMPP is employed to model video sources with scene changes. Simulation studies to analyze statistical multiplexing using the first and second order AR models are in [Nomura 89a] and [Ogino 88]. The advantages of AR models are easy model identification and simulation.



## 6.2 Models of superposed processes and general burst traffic

To evaluate statistical multiplexing, modeling of superposed sources is necessary. In addition, superposition of heterogeneous traffic must be modeled. There are two alternatives: superpose at the level of statistics, such as, average and variance, and model after that; and model individual sources, and superpose at the level of models.

An example of the former is a renewal approximation [Whitt 83]. Examples of the latter are a PH-Markov renewal process [Machihara 88] and a Markovian arrival process [Lucantoni 88], which include MMPPs, and a discrete-time model. They have the characteristic that their superpositions are contained in them.

Examples of MMPPs which model heterogeneous traffic are provided by [Arai 89] and [Saito 91]. In particular, performance measures of individual classes are presented by [Saito 91].

Discrete-time models for ATM nodes are employed by [Hirano 89] and [Murata 89,90]. General discrete-time models which include these are provided by [Ushijima 72] and [Morris 81]. The discrete-time models have the advantage that they can easily reflect the fact that a cell is of fixed length.

## 7. Conclusions

This chapter addressed the issues of transmitting voice and data in public integrated networks by circuit switching, hybrid switching, burst switching and packet switching. The emphasis was on reviewing performance considerations for integrated networks.

In comparison to dedicated networks, service and network integration have major advantages in economic planning, development, implementation and maintenance. The introduction of fiber optics has permitted new services employing wide bandwidths.

After achieving narrowband integration, interest is moving toward broadband integration. In broadband ISDN (BISDN), video services as well as voice and data traffic can be made available commercially. Modelling and performance evaluation on integrated networks including video traffic becomes important. We also discussed it.

ATM is now considered the most promising approach for integration of future broadband telecommunication networks. However, there are some design and performance problems to be solved. Current interests focus on the performance evaluation of multi-media traffic in ATM networks.



## Appendix 1. Speech activity model

In normal telephone conversation, the average talkspurt activity of a single user is well known to be less than half. There are several telecommunications techniques (e.g. TASI [Bullington 59], hybrid switching with SAD (see Section 3), burst switching (see Section 4) and packet switching (see Section 5) that utilize this property with speech activity detection (SAD) to raise the utilization of transmission links.

Therefore, for dimensioning telecommunications resources and performance evaluation of telecommunications methods using SAD, the stochastic properties of a single voice source should be investigated and a model for generating on-off speech patterns must be developed (Table 4).

In his pioneering work [Brady 65,68,69] Brady used a fixed threshold speech detector with 15 ms throwaway (all spurts less than 15 ms are presumed to be noise and are rejected) and 200 msec fillin (all silence less than 200 ms is filled in). He showed that the simplest Markovian model for an on-off pattern with four states depending on who is talking is not generally supported by the data [Brady 68]. The exponential distribution is an empirical fit to talkspurts, but not to silence period, and a six-state model is proposed [Brady 69]. Although these results are elaborately examined, we should note that the analyzed results are sensitive to the design of the speech detector [Brady 65].

In [Yatsuzuka 82a], a statistical analysis as well as a proposal for a speech detector is shown. The important feature of the analysis is that although the talkspurt coefficient of variation is slightly less than one, the coefficient of variation of the silence distribution is almost three. This result was the motivation of [O'Reilly 87a] (see Section 4). [Gruber 82]

also identified the distributions of talkspurt and silence durations with zero hangover and fillin, which enabled parameters to be computed with any value of hangover and fillin. His conclusions for the case of zero hangover and fillin are that the measured talkspurt pdf can be modeled approximately by a geometric pdf, and that the measured pdf for silent periods can be modeled approximately by two suitably weighted geometric pdfs.

A similar result was obtained by Lee et al. [Lee 86]. They used 14 two-way telephone conversations in English and SAD with no fillin and no hangover. They reported that the pdfs of talkspurt and silence durations can be expressed by two weighted geometric functions, when no fillin and no hangover SAD is used. For SAD with larger hangover (fillin) greater than 200 ms, silence (talkspurt) duration can be considered to be exponentially distributed and the pdf of talkspurt (silence) duration can be modeled by a constant-plus-exponential.

Lastly, the simulation voice model based on an unpublished experimental study [Seguel 82] is discussed here (see Section 5). Seguel et al. used a four-state model (a talkspurt, a silence, a pause or an activity) instead of a two-state model (a talkspurt or a silence period, see Figure 4) for a single voice source. They define a talkspurt as a time period in which a person is speaking, a silence as a time period in which a person is not speaking, a pause as the time period of a gap due to stop consonants and slight hesitation during speech, and an activity as a time period of unbroken voice. The probability distribution of these variables was taken from an earlier unpublished study of theirs: the length of a talkspurt has an exponential distribution with the mean value = 4 s, the length of a silence has the same distribution, the length of an activity has an Erlang distribution with a phase



of 3 and a mean value of 225 ms, and the length of a pause has an exponential distribution with a mean of 60 ms.

The characteristics of on-off patterns in telephone conversation depend on the language spoken [Lyghounis 74], the sex of speakers [Brady 68], and the design of a speech activity detector. They are quite delicate. However, most studies agree on the hyperexponential modeling of talkspurt and silence durations.

## Appendix 2. Flow control with bit reduction

One of the essential differences between voice communication and data communication lies in the fact that user information itself in voice communication is redundant.

Hence, it is natural to consider the flow control reducing the carried information of voice traffic when the network is congested. In integrated networks including voice and/or video, flow control utilizing this property is an attractive technique.

The length of a voice packet or the length of a voice packet generation interval is variable, or a low priority voice packet is discarded in a system using such a control [Webber 77, Dubnowski 79, Bially 80b, Cox 80, Goodman 80, Seguel 82, Sato 88]. In DSI, bit reduction can generate the capacity to transmit talkspurts which would otherwise be clipped [Lyghounis 74, Weinstein 80, Yatsuzuka 82b, Kou 85].

Analyses are given in [Campanella 76, Bially 80a, Jayant 81, Holtzman 85, Fredericks 86, Yin 87a,b, Sriram 88] and [Saito 89a,90b]. For DSI, Campanella indicates that the occurrence probability of clips longer than 50 ms drops significantly when bit reduction is used. Embedded coding, the most practical method of bit reduction, of the packet voice system is analyzed in [Bially 80a, Yin 87b, Sriram 88], and [Saito 89a,90b] as well as in an

experimental work on WB SATNET [Weinstein 83]. Bially et al. evaluate the control using embedded coding by simulation. The simulated network has a star configuration. Yin et al. investigate the performance of the flow control using embedded coding, odd-even samples or speech energy detection based on the measurement of the speech activity (the number of active talkers) or the buffer content [Yin 87b] at a packet voice multiplexer. In this analysis, low priority packets are assumed to be discarded during overload. The analysis uses the fluid flow approximation employed in [Anick 82], and [Gaver 82]. Speech energy detection with speaker activity measurement is superior in light traffic, and embedded coding or even-odd samples with buffer content measurement is superior in heavy traffic. Control using embedded coding with buffer content measurement is also analyzed in [Sriram 88]. Here, however, the packet length is assumed to be reduced during overload based on the measurement just before the packet transmission. Sriram suggested the possibility of approximating the aggregated arrival process by a Poisson process when bit dropping is used. The multiplexer is modeled by using an M/D/1/K model in which the service rate is state-dependent. An optimal control using embedded coding is given by [Saito 89a,90b]. He investigates an optimization problem with the constraint

$$\max E[\text{coding rate}], \text{ subj. to } E[\text{delay}] \leq T$$

and shows that the optimal control has a simple structure, linear feedback with saturation.

There are several related works. The voice quality under the flow control which reduces the coding rate of a voice packet when a queue length exceeds a specified threshold is evaluated in [Holtzman 85]. The multiplexer is modeled by using an M/M/1 queueing model in which the arrival and service rates are state-dependent in this analysis. [Fredericks 86]



examined the flow control which provides the coding rate by

$$\min [\text{coding rate during normal load}, K/\text{the number of active talkers}],$$

where  $K$  is the maximum bandwidth. He derives the mean sojourn time and the fraction of time during which the voice coding rate is below the admissible level. A delay scheme of low priority packets during overload is given in [Yin 87a]. This analysis uses a fluid flow approximation.

Experimental studies are reported in [Ueda 83], and [Muise 86]. Speech energy detection is tried in [Ueda 83]. [Muise 86] reports that bit dropping has the effect that the load/service curves do not show typical cliff behavior, and that bit dropping provides graceful degradation as the load is increased.

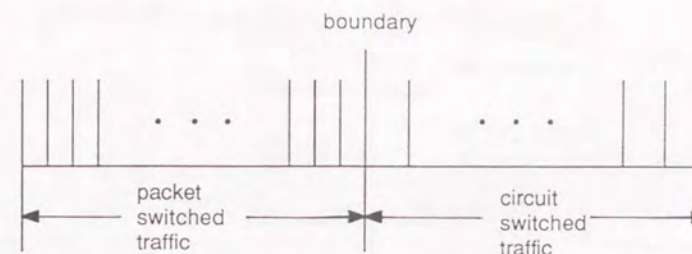


Figure 1. Frame structure

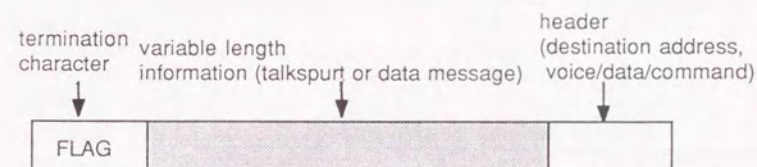


Figure 2. A burst



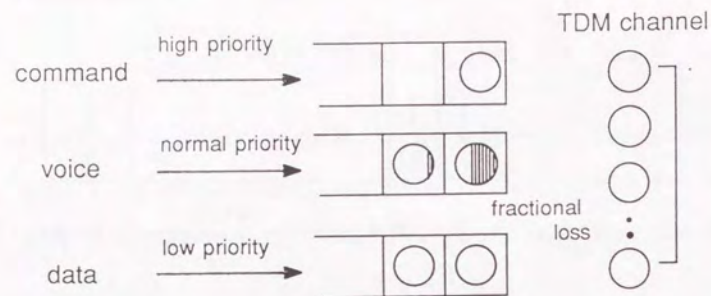


Figure 3. Burst switching

Table 1. Performance evaluation of burst switching

	single node		network
	with voice delay	without voice delay	
exact	Descloux 85	Lim 86 Aboul-Magd 88	
approximation	Ma 87	(Lehoczký 81b) (Sriram 83) (Williams 84) (Li 85) (O'Reilly 85, 86b) O'Reilly 87a (McDaysan 88) Aboul-Magd 88 Ma 88	Leon-Garcia 86
comparative studies between burst and fast packet switching (approximation)		Li 87 O'Reilly 87b	
simulation			Morse 85 Jack 86

( ): See "hybrid switching with SAD".



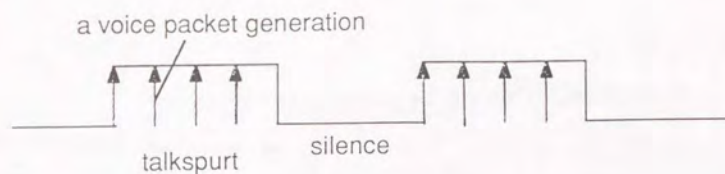


Figure 4. Voice source behavior

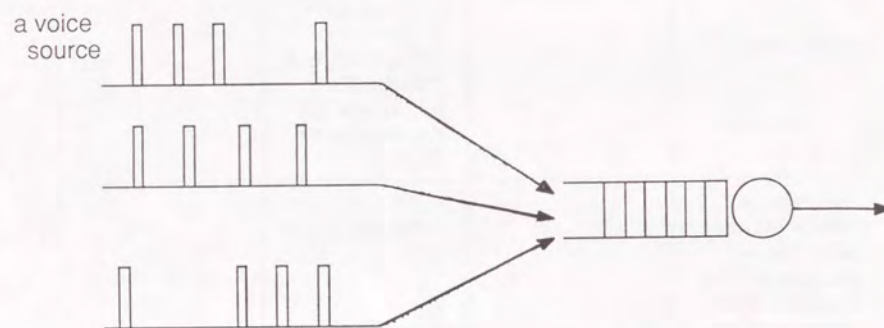


Figure 5. A statistical packetized voice mutiplexer

Table 2. Performance evaluation of a packet voice multiplexer

	model	solution technique
Tanaka 82	Discrete time model	numerical evaluation of the state equations
Stern 83,84	Semi-Markov model	functional iteration /spectral factorization
Jenq 84	M/G/1 Discrete time model GI/G/1	QNA
Daigle 85, 86	Semi-Markov model MMPP/M/1 Fluid flow model	matrix geometry
Sriram 86	GI/G/1	QNA
Heffes 86	SPP/G/1	matrix geometry
Ide 88	N-IPP/G/1	supplementary variable method
Li 88d	Fluid flow / discrete time finite buffer 2-dimensional Markov chain	

SPP: Switched Poisson Process = two-state MMPP

IPP: Interrupted Poisson Process

N : Number of off-hook users



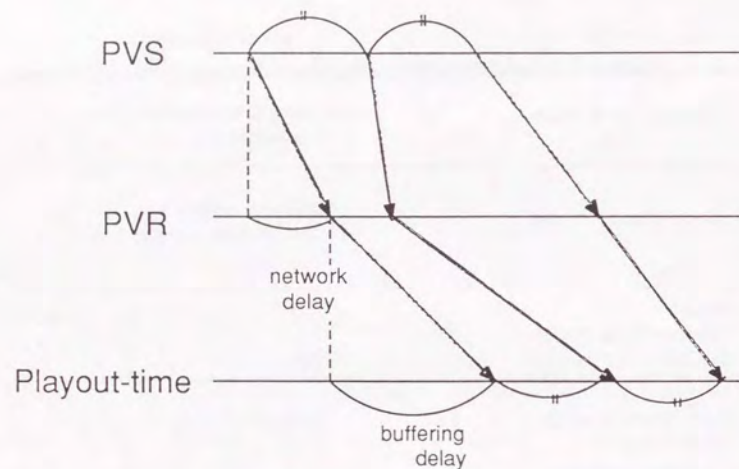


Figure 6. Timing diagram of playout

Table 3. Video source model

		Bit rate statistics (Mbits/s)			Model
		average	standard deviation	max	
Maglaris [88]	Sen [89]	3.9	1.7	10.6	first order AR 2 dimensional Markovian Process
Verbiest [88] (TV conference) (CATV)		4.2		14	
	(video phone)	16.8	4.3	44.7	
[89] (video conf.)		4.3	2.9	19.1	
	(Studio TV)	4.9	2.3	15.1	
		26.5	3.4	51.5	
Ogino [88]		0.246	0.076	0.441	2nd order AR
	[87]	0.062	0.038	0.236	first-order AR
		0.139	0.035	0.265	
		0.253	0.077	0.441	
Nomura [89a]		13.3	4.1		
		10.3	4.0		
	[89b]	2.7	1.08	9.2	
		1.7	0.34	5.3	
		2.0	0.8	5.0	
		1.2	0.36	3.12	
		0.25	0.075	0.45	
		0.14	0.042	0.266	



Figure 7. Video source model



Table 4. On-off pattern in speech

	Brady	Yatsuzuka	Gruber	Lee
talkspurt	exponential	c.v. $\leq 1$	geometric	two weighted geometric
silence		c.v. $\approx 3$	two weighted geometric	two weighted geometric
SAD	throwaway=15ms fillin=200ms fixed threshold	hangover=32ms short time energy zero-crossing rate sign bit sequence	hangover=0 fillin=0 fixed threshold	hangover=0 fillin=0 short time energy zero-crossing rate

## Chapter III

### The Departure Process of an N/G/1 Queue

*The departure process of an N/G/1 queue is investigated. The arrival process called an N process is a versatile point process and includes, for example, a Markov-modulated Poisson process, which is a model of a voice packet arrival process. The first passage analysis yields LSTs of distributions of the interdeparture times and counts of departure. Emphasis is on the interdeparture times of an N/D/1 queue. Numerical examples show that correlation of interarrival times is likely to be preserved in interdeparture times. [Saito 89b]*



## 1. Introduction

Study of departure is important, because in a network of queueing systems, the arrival at a particular queue in the network is the departure of one or more other queueing systems plus traffic from external sources [Whitt 83]. Analysis of departure processes, stimulated by the pioneering work of Burke [Burke 56], has yielded many useful results. For a review of related literature, see [Daley 76] and [Disney 85,87].

This paper addresses the departure processes from  $N/G/1$  queues. Emphasis is on the interdeparture times of an  $N/D/1$  queue, which is directly applicable to the analysis of an ATM network [Saito 91] as well as to the departure counts. The arrival process, called an  $N$ -process, is a versatile point process recently introduced by Neuts [Neuts 79]. The  $N$ -process contains, for example, the Markov modulated Poisson processes (MMPP) as a special case, which is useful in modeling arrival processes in communications networks. The interrupted Poisson process in a telephone engineering model is an example of an MMPP [Kuczura 73]. Data traffic is also modeled as an MMPP [Heffes 80, Rossiter 87]. In particular, recently, MMPPs have been employed as models of the arrival process for voice and video cells ([Daigle 86, Heffes 86, Ide 88, Yamada 89a, Saito 91]). Thus, the importance of the analysis of an  $N$ -process has been increasing.

Analysis is done by the first passage time technique. The process is characterized by  $C_P(z)$ , a index based on the covariance of interevent times.

## 2. Model and notation

An  $N$ -process is defined based on [Neuts 79] and [Ramaswami 80]. Consider a continuous-time Markov process, called a phase process, with state space  $\{1, \dots, m, m+1\}$  for which the states  $\{1, \dots, m\}$  are transient and the state  $m+1$  is absorbing. The infinite generator of such a Markov process then has the form

$$\begin{pmatrix} T & \mathbf{T}^0 \\ 0 & 0 \end{pmatrix}. \quad (2.1)$$

$T$  corresponds a transition between transient states, and is an  $(m, m)$  matrix with  $T_{ii} < 0$  and  $T_{ij} \geq 0$  for  $i \neq j$  such that  $T^{-1}$  exists. The vector  $\mathbf{T}^0$  is non negative and satisfies

$$T\mathbf{e} + \mathbf{T}^0 = 0, \quad (2.2)$$

where  $\mathbf{e} = (1, \dots, 1)'$ .  $\mathbf{T}^0$  corresponds to a transition from transient states to the absorbing state. A vector  $(\mathbf{a}, a_{m+1})$  of initial probability is also given and satisfies

$$\mathbf{a}\mathbf{e} + a_{m+1} = 1. \quad (2.3)$$

In what follows, we assume  $a_{m+1} = 0$ . During any sojourn in the transient state  $i$ ,  $1 \leq i \leq m$ , there are Poisson process arrivals of rate  $\lambda_i$  and group size density  $\{p_i(k); k \geq 0\}$ . Let  $\phi_i(z)$  denote the p.g.f. of  $\{p_i(k)\}$  and let  $\phi = \{\phi_1(z), \dots, \phi_m(z)\}$ . At  $(i, j)$ -renewal transitions (that is, from the transient state  $i$  to the 'instantaneous' state  $m+1$ , and from the state  $m+1$  to the transient state  $j$ ), there are group arrivals with density  $\{\tau_{ij}(k); k \geq 0\}$  where p.g.f. is  $\Phi_{ij}(z)$ . Let  $\Phi(z)$  denote the  $(m, m)$  matrix of entries  $\Phi_{ij}(z)$ . At  $(i, j)$ -transitions,  $i \neq j$ , there are group arrivals with size density  $\{q_{ij}(k); k \geq 0\}$  where p.g.f. is  $\psi_{ij}(z)$ . For notational convenience in the sequel, we set  $\psi_{ii} \equiv 1$ ,  $1 \leq i \leq m$  and let  $\psi(z)$  denote the  $(m, m)$  matrix of entries  $\psi_{ij}$ .



Let  $N(t)$  and  $J(t)$ ,  $t \geq 0$ , denote respectively the number of arrivals in  $(0, t]$  and the phase at  $t$ . Then,  $(N(t), J(t))$  is a Markov process and the generating function for the number of arrivals is

$$\tilde{P}(z, t) \triangleq \sum_{k=0}^{\infty} z^k P(k, t)$$

where  $P(k, t) \triangleq (p_{ij}(k, t))$ ,  $p_{ij}(k, t) \triangleq \Pr(N(t) = k, J(t) = j | N(0) = 0, J(0) = i)$  is given by

$$\tilde{P}(z, t) = \exp(R(z)t) \quad |z| \leq 1 \quad (2.4)$$

with  $R(z) = \Delta(\lambda)\Delta(\phi(z)) - \Delta(\lambda) + T \circ \psi(z) + T^0 A^0 \circ \Phi(z)$  where  $\Delta(\lambda) = \text{diag}(\lambda_1, \dots, \lambda_m)$ ,  $\Delta(\phi(z)) = \text{diag}(\phi_1(z), \dots, \phi_m(z))$ ,  $A^0 = \text{diag}(a_1, \dots, a_m)$  and  $T^0 = (T^0, \dots, T^0)$  [Neuts 79], [Ramaswami 80]. Here,  $\circ$  denotes the entrywise product of two matrices.

From now we assume non-triviality of the  $N$ -process, viz,  $R(0) \neq Q^*$ , where

$$Q^* = T + T^0 A^0.$$

Then, we can obtain a useful result of this assumption.

**Collary 1.3.18** [Ramaswami 80].  $(sI - R(0))^{-1}$  exists for all  $s \geq 0$ .

These arrivals according to an  $N$ -process join an FIFO single server queue with  $K$  waiting rooms,  $K \leq \infty$ . If a whole batch cannot be accepted because its size is larger than the number of unoccupied waiting rooms, the batch fills idle servers and unoccupied waiting rooms, and remaining customers in the batch are rejected and lost (PBAS).

The service times of all customers are independent and identically distributed with distribution function  $\tilde{H}(\cdot)$ . The Laplace-Stieltjes transform (LST) of  $\tilde{H}(\cdot)$  is denoted by  $\tilde{H}(\cdot)$ . Let  $h$  be the mean service time.

We refer the departure process of this queueing system as the sequence of the service completion epoch. That is, the arrivals which cannot find unoccupied waiting rooms and leave the system immediately, are not counted.

### 3. Preliminaries

We define  $\{\tau_n : n \geq 0\}$  as the successive epochs of departure (with  $\tau_0 = 0$ ). We further define  $X_n$  and  $J_n$  to be, respectively, the number of customers in the system and the phase of the  $N$ -process just after  $\tau_n$ . Set  $t_n \triangleq \tau_{n+1} - \tau_n$ . Then, the sequence  $\{(X_n, J_n, t_n) : n \geq 0\}$  forms a semi-Markov sequence and transition probability matrix  $\tilde{Q}(\cdot)$  given by

$$\tilde{Q}(z) = \begin{pmatrix} \tilde{B}_0(z) & \tilde{B}_1(z) & \tilde{B}_2(z) & \cdots & \sum_{k=K}^{\infty} \tilde{B}_k(z) \\ \tilde{A}_0(z) & \tilde{A}_1(z) & \tilde{A}_2(z) & \cdots & \sum_{k=K}^{\infty} \tilde{A}_k(z) \\ 0 & \tilde{A}_0(z) & \tilde{A}_1(z) & \cdots & \sum_{k=K-1}^{\infty} \tilde{A}_k(z) \\ \vdots & & \ddots & & \\ 0 & 0 & \cdots & \tilde{A}_0(z) & \sum_{k=1}^{\infty} \tilde{A}_k(z) \end{pmatrix}, z \geq 0, \quad (3.1)$$

where the  $(m, m)$  matrices of mass functions are

$$\tilde{A}_n(z) = \int_0^z P(n, u) d\tilde{H}(u) \quad n \geq 0, z \geq 0 \quad (3.2)$$

$$\tilde{B}_n(z) = \sum_{k=1}^{n+1} (\tilde{U}_k \star \tilde{A}_{n-k+1})(z) \quad n \geq 0, z \geq 0 \quad (3.3)$$

and

$$\tilde{U}_k(z) = \left\{ \int_0^z P(0, y) dy \right\} \{T^0 A^0 \circ r(k) + T \circ q(k) + \Delta(\lambda)\Delta(p(k))\} \quad k \geq 1, z \geq 0 \quad (3.4)$$

where  $r(k)$  and  $q(k)$  are  $(m, m)$  matrices with respective entries  $r_{ij}(k)$  and  $q_{ij}(k)$ ,  $p(k)$  is an  $m$ -vector with entries  $p_i(k)$  and  $\Delta(p(k))$  is an  $(m, m)$  diagonal matrix with  $p(k)$  along



the diagonal. Also  $\star$  in the definition of  $\tilde{B}_n(\mathbf{z})$  denotes matrix convolution. Let  $Q(\cdot)$ ,  $A_n(\cdot)$ ,  $B_n(\cdot)$  and  $U_k(\cdot)$  denote the LSTs of  $\tilde{Q}(\cdot)$ ,  $\tilde{A}_n(\cdot)$ ,  $\tilde{B}_n(\cdot)$  and  $\tilde{U}_k(\cdot)$ .

Before the analysis of the departure process of an  $N/G/1$  queue, we derive results on the first passages of the semi-Markov process  $\tilde{Q}(\cdot)$ . Consider the process  $(X_n, J_n)$  and define level  $\mathbf{i}$  as the set of the states  $\mathbf{i} \triangleq \{(i, j); 1 \leq j \leq m\}$ , where  $0 \leq i \leq K+1$ . Let  $\tilde{G}_{j,j'}^{[i,i']}(k, \mathbf{z})$  be the probability that, given that the semi-Markov process  $\tilde{Q}(\cdot)$  starts in the state  $(i, j)$ , it reaches the level  $\mathbf{i}'$  for the first time after  $k$  transitions by visiting  $(i', j')$  and the time of such a first passage is at most  $\mathbf{z}$ . The matrix  $\tilde{G}^{[i,i']}(k, \mathbf{z})$  has entries  $\tilde{G}_{j,j'}^{[i,i']}(k, \mathbf{z})$ ,  $1 \leq j, j' \leq m$  and is the probability that the number of transitions of the first passage from level  $\mathbf{i}$  to level  $\mathbf{i}'$  is  $k$  and its time is at most  $\mathbf{z}$ . We define the transform as

$$G^{[i,i']}(z, s) = \sum_{k=1}^{\infty} \int_0^{\infty} e^{-sz} d\tilde{G}^{[i,i']}(k, \mathbf{z}) z^k,$$

for  $|z| \leq 1$  and  $\Re s \geq 0$ .

To simplify the notation, we set

$$G^{[i,i]}(z, s) = I.$$

Thus, ordinary arguments on the first passage time yield the following equations.

$$G^{[i,i']}(z, s) = G^{[i,i-1]}(z, s) \cdots G^{[i'+1,i']}(z, s) \quad 0 \leq i' < i \leq K \quad (3.5)$$

$$G^{[i,i-1]}(z, s) = z \sum_{k=0}^{K-i} A_k(s) G^{[i+k-1,i-1]}(z, s) + z \sum_{k=K-i+1}^{\infty} A_k(s) G^{[K,i-1]}(z, s) \quad (3.6)$$

When  $K = \infty$ , Eq. (3.6) is consistent with the result in [Ramaswami 80].

Let  $\tilde{W}_{j,j'}(k, \mathbf{z})$  be the probability that, given that the semi-Markov process  $\tilde{Q}(\cdot)$  starts in the state  $(0, j)$ , it is in the state  $(0, j')$  after  $k$  transitions and the time is at most  $\mathbf{z}$ .

The matrix  $\tilde{W}(k, \mathbf{z})$  has entries  $\tilde{W}_{j,j'}(k, \mathbf{z})$ ,  $1 \leq j, j' \leq m$ . Define the transform as

$$W(z, s) = \sum_{k=0}^{\infty} \int_0^{\infty} e^{-sz} d\tilde{W}(k, \mathbf{z}) z^k,$$

for  $|z| \leq 1$ , and  $\Re s \geq 0$ .

$W(z, s)$  satisfies

$$W(z, s) = I + z \left( \sum_{k=0}^{K-1} B_k(s) G^{[k,0]}(z, s) + \sum_{k=K}^{\infty} B_k(s) G^{[K,0]}(z, s) \right) W(z, s). \quad (3.7)$$

Substitute Eq. (3.3) into Eq. (3.7), and use Eqs. (3.5)-(3.6). Consequently,

$$W(z, s) = \left\{ I - \left( \sum_{k=0}^{K-1} U_{k+1}(s) G^{[k+1,0]}(z, s) + z \sum_{k=K}^{\infty} U_{k+1}(s) \left( \sum_{l=0}^{\infty} A_l(s) \right) G^{[K,0]}(z, s) \right) \right\}^{-1}. \quad (3.8)$$

#### 4. The interdeparture times of an $N/G/1$ queue

We focus on the stationary probability of interdeparture times. We assume in the remaining of this paper that the Markov process  $(X_n, J_n)$  is stationary. For infinite waiting rooms it is supposed that the traffic intensity  $\rho \triangleq \theta R'(1)eh < 1$ , where  $\theta$  is the invariant vector of the Markov Process  $Q^*$  [Ramaswami 80]. The stationary density of the number of customers in the system just after the departure is denoted by  $\mathbf{x} = (x_0, x_1, \dots, x_K)$ .

$$\mathbf{x}Q(0) = \mathbf{x}, \quad \mathbf{x}e = 1 \quad (4.1)$$

Here,  $\mathbf{x}_k$  is the stationary probability vector with entries  $x_k(j)$ ,  $j = 1, \dots, m$ , where the number of customers in the system is  $k$  and the phase is  $j$ .

Let  $\tilde{D}(k, t)$  be the stationary distribution that the sum of the length of  $k$ -consecutive interdeparture intervals,  $t_1 + \dots + t_k$ , is at most  $t$ . We define

$$D(z, s) \triangleq \sum_{k=1}^{\infty} \int_0^{\infty} e^{-st} d\tilde{D}(k, t) z^k.$$



Employing the stationary distribution  $\mathbf{x}$  and  $Q(s)$  defined in the previous section, we obtain,

$$D(z, s) = \mathbf{x} \sum_{k=1}^{\infty} Q^k(s) z^k \mathbf{e}. \quad (4.2)$$

Using the facts that

$$Q(s)\mathbf{e} = \begin{pmatrix} \sum_{k=0}^{\infty} A_k(s)\mathbf{e} \\ \sum_{k=0}^{\infty} A_k(s)\mathbf{e} \\ \vdots \\ \sum_{k=0}^{\infty} A_k(s)\mathbf{e} \end{pmatrix} + \begin{pmatrix} S(s) \sum_{k=0}^{\infty} A_k(s)\mathbf{e} \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad (4.3)$$

and

$$\sum_{k=0}^{\infty} A_k(s)\mathbf{e} = H(s)\mathbf{e} \quad (4.4)$$

where

$$\begin{aligned} S(s) &= (sI - R(0))^{-1}(R(1) - R(0)) - I \\ &= -(sI - R(0))^{-1}(sI - R(1)), \end{aligned} \quad (4.5)$$

we obtain

$$Q^n(s)\mathbf{e} = H^n(s)\mathbf{e} + \sum_{k=0}^{n-1} H^{k+1}(s)Q^{n-k-1}(s) \begin{pmatrix} S(s)\mathbf{e} \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \quad n \geq 1. \quad (4.6)$$

We note here that  $R(1)\mathbf{e} = 0$ . Thus,

$$S(s)\mathbf{e} = -s(sI - R(0))^{-1}\mathbf{e}. \quad (4.7)$$

By Eq. (4.6),

$$D(z, s)$$

$$\begin{aligned} &= \mathbf{x} \sum_{n=1}^{\infty} H^n(s) z^n \mathbf{e} + z H(s) \mathbf{x} \sum_{n=1}^{\infty} \sum_{k=0}^{n-1} H^k(s) Q^{n-1-k}(s) z^{n-1} \begin{pmatrix} S(s)\mathbf{e} \\ 0 \\ \vdots \\ 0 \end{pmatrix} \\ &= \frac{z H(s)}{1 - z H(s)} + \frac{z H(s)}{1 - z H(s)} \mathbf{x} \sum_{k=0}^{\infty} Q^k(s) z^k \begin{pmatrix} S(s)\mathbf{e} \\ 0 \\ \vdots \\ 0 \end{pmatrix} \\ &= \frac{z H(s)}{1 - z H(s)} + \frac{z H(s)}{1 - z H(s)} \sum_{k=0}^K \mathbf{x}_k G^{[k,0]}(z, s) W(z, s) S(s) \mathbf{e}. \end{aligned} \quad (4.8)$$

In particular, the first order of  $z$  in  $D(z, s)$  which provides the LST of the stationary distribution of interdeparture times, is derived from Eq (4.8) as,

$$\begin{aligned} D(z, s) &= z H(s) \{1 + \mathbf{x}_0 S(s)\mathbf{e}\} + O(z^2) \\ &= z H(s) \left\{ \sum_{k=1}^K \mathbf{x}_k \mathbf{e} + \mathbf{x}_0 (sI - R(0))^{-1} (R(1) - R(0)) \mathbf{e} \right\} + O(z^2). \end{aligned} \quad (4.9)$$

The first term on the right-hand side of Eq. (4.9) denotes the event that there are customers in the system just after departure and the second term denotes that there is no customer in the system just after departure and that an idle period starts. Equation (4.9) is an intuitive result.

We can obtain the mean interdeparture time,  $\bar{d}$ , by the first order of  $z$  in Eq. (4.9) from Eq. (4.7).

$$\begin{aligned} \bar{d} &= -\frac{\partial}{\partial s} [H(s) \{1 + \mathbf{x}_0 S(s)\mathbf{e}\}]_{s=0} \\ &= h - \mathbf{x}_0 R^{-1}(0)\mathbf{e}. \end{aligned} \quad (4.10)$$

$\bar{d}$  also means the mean interarrival time of accepted customers as well as the mean interdeparture time. For  $K = \infty$ , Eq. (4.10) is consistent with the result in [Ramaswami 80].



Similarly, the variance of interdeparture times,  $var$ , is given by,

$$var = H''(0) - 2h\mathbf{x}_0 R^{-1}(0)e + 2\mathbf{x}_0 R^{-2}(0)e - \bar{d}^2. \quad (4.11)$$

## 5. The interdeparture times of an N/D/1 queue

This section focuses on a deterministic service case, which is important in applications [Saito 91]. In this case,

$$H(s) = e^{-sh}. \quad (5.1)$$

In a simplified notation,

$$G^{[i,j]}(z) \triangleq G^{[i,j]}(z, 0), \quad W(z) \triangleq W(z, 0). \quad (5.2)$$

We first note that

$$\frac{\mathbf{x}_0}{1-z} = \sum_{k=0}^K \mathbf{x}_k G^{[k,0]}(z) W(z). \quad (5.3)$$

Then, using the fact that during the first passage from set  $i$  to set  $0$ , the server is not idle and that the time for serving  $n$  customers is  $nh$ ,

$$G^{[i,j]}(z, s) = G^{[i,j]}(ze^{-sh}). \quad (5.4)$$

Therefore, note that

$$\sum_{l=0}^{\infty} A_l(s) = e^{-sh} e^{R(1)h}, \quad (5.5)$$

and use Eqs. (3.8), (5.3) and (5.4),

$$\begin{aligned} & \sum_{k=0}^K \mathbf{x}_k G^{[k,0]}(z, s) \\ &= \sum_{k=0}^K \mathbf{x}_k G^{[k,0]}(ze^{-sh}) \end{aligned}$$

$$\begin{aligned} &= \frac{\mathbf{x}_0}{1-ze^{-sh}} W^{-1}(ze^{-sh}) \\ &= \frac{\mathbf{x}_0}{1-ze^{-sh}} W^{-1}(ze^{-sh}, 0) \\ &= \frac{\mathbf{x}_0}{1-ze^{-sh}} \left\{ I - \left( \sum_{l=0}^{K-1} U_{l+1}(0) G^{[l+1,0]}(ze^{-sh}) + ze^{-sh} \sum_{l=K}^{\infty} U_{l+1}(0) e^{R(1)h} G^{[K,0]}(ze^{-sh}) \right) \right\} \\ &= \frac{\mathbf{x}_0}{1-ze^{-sh}} (I - L(ze^{-sh})), \end{aligned} \quad (5.6)$$

where

$$L(z) = \sum_{l=0}^{K-1} U_{l+1}(0) G^{[l+1,0]}(z) + z \sum_{l=K}^{\infty} U_{l+1}(0) e^{R(1)h} G^{[K,0]}(z). \quad (5.7)$$

$L(\cdot)$  denotes the generating function of the number of served customers in a busy period.

Consequently, substituting Eqs. (5.6) and (4.7) into (4.8), substituting (5.5) into (3.8), using (5.4) in (3.8), and employing the fact that

$$U_k(s) = -(sI - R(0))^{-1} R(0) U_k(0),$$

we obtain

$$\begin{aligned} D(z, s) &= \frac{ze^{-sh}}{1-ze^{-sh}} \\ &- s \frac{ze^{-sh}}{(1-ze^{-sh})^2} \mathbf{x}_0 (I - L(ze^{-sh})) \{sI - R(0)(I - L(ze^{-sh}))\}^{-1} e. \end{aligned} \quad (5.8)$$

$L(1)$  is stochastic and the mean number of transitions in the semi-Markov process before the first return to states in level  $0$ , starting from level  $0$  is given by  $[\frac{\partial}{\partial z} L(z)]_{z=1}$ . Thus, the vector  $\mathbf{x}_0$  can be determined by the equations [Ramaswami 80],

$$\mathbf{x}_0 L(1) = \mathbf{x}_0 \quad (5.9)$$

and

$$\mathbf{x}_0 \left[ \frac{\partial}{\partial z} L(z) \right]_{z=1} e = 1. \quad (5.10)$$



## 5.1 Mean interdeparture time

Employing Eq. (5.8),

$$\begin{aligned} & \sum_{k=1}^{\infty} E[t_1 + \dots + t_k] z^k \\ &= - \left[ \frac{\partial}{\partial s} D(z, s) \right]_{s=0} \\ &= \frac{zh}{(1-z)^2} - \frac{z}{(1-z)^2} \mathbf{x}_0 R^{-1}(0) \mathbf{e}. \end{aligned} \quad (5.11)$$

Here, we note that, based on the assumption of stationarity,

$$E[t_1 + \dots + t_k] = k\bar{d}.$$

Thus Eq. (5.11) is consistent with Eq. (4.10).

## 5.2 Second moment of the length of the interdeparture times

Using Eq. (5.8),

$$\begin{aligned} & \sum_{k=1}^{\infty} E[(t_1 + \dots + t_k)^2] z^k \\ &= \left[ \frac{\partial^2}{\partial s^2} D(z, s) \right]_{s=0} \\ &= \frac{h^2 z(1+z)}{(1-z)^3} - 2h \frac{(1+z)z}{(1-z)^3} \mathbf{x}_0 R^{-1}(0) \mathbf{e} \\ & \quad + \frac{2z}{(1-z)^2} \mathbf{x}_0 R^{-1}(0) (I - L(z))^{-1} R^{-1}(0) \mathbf{e} \\ &= z \{ h^2 - 2h \mathbf{x}_0 R^{-1}(0) \mathbf{e} + 2 \mathbf{x}_0 R^{-2}(0) \mathbf{e} \} \\ & \quad + z^2 \{ 4h^2 - 8h \mathbf{x}_0 R^{-1}(0) \mathbf{e} + 4 \mathbf{x}_0 R^{-2}(0) \mathbf{e} + 2 \mathbf{x}_0 R^{-1}(0) U_1(0) A_0(0) R^{-1}(0) \mathbf{e} \} \\ & \quad + O(z^3) \end{aligned} \quad (5.12)$$

Here, we use the fact that

$$\begin{aligned} L(z) &= U_1(0) G^{[1,0]}(z) + O(z^2) \\ &= z U_1(0) A_0(0) + O(z^2), \end{aligned} \quad (5.14)$$

where  $A_0(s) = e^{-sh} e^{R(0)h}$ . The first order of  $z$  in Eq. (5.13) provides the variance of interdeparture times,

$$var = 2 \mathbf{x}_0 R^{-2}(0) \mathbf{e} - (\mathbf{x}_0 R^{-1}(0) \mathbf{e})^2, \quad (5.15)$$

which is consistent with Eq. (5.13) with  $H(s) = e^{-sh}$ .

Employing Eq. (5.12), characteristics of the length of interdeparture time can be derived. For example, let  $c_k \triangleq \text{cov}(t_1, t_{k+1})$  be the covariance with lag  $k$  of the length of interdeparture times. Then using that

$$\left[ \frac{\partial^2}{\partial s^2} D(z, s) \right]_{s=0} = \frac{z}{(1-z)^2} var + \frac{2z}{(1-z)^2} \sum_{k=1}^{\infty} c_k z^k + \frac{z(1+z)}{(1-z)^3} \bar{d}^2, \quad (5.16)$$

Eq. (5.12) yields

$$\begin{aligned} C_P(z) &\triangleq \frac{\sum_{k=0}^{\infty} c_k z^k}{\bar{d}^2} \\ &= \frac{1}{\bar{d}^2} \left\{ \frac{(1-z)^2}{2z} \left[ \frac{\partial^2}{\partial s^2} D(z, s) \right]_{s=0} + \frac{1}{2} var - \frac{1+z}{2(1-z)} \bar{d}^2 \right\} \\ &= \frac{1}{\bar{d}^2} \left\{ \mathbf{x}_0 R^{-1}(0) (I - L(z))^{-1} R^{-1}(0) \mathbf{e} - \frac{(\mathbf{x}_0 R^{-1}(0) \mathbf{e})^2}{1-z} + \mathbf{x}_0 R^{-2}(0) \mathbf{e} \right\} \end{aligned} \quad (5.17)$$

For a Poisson process,  $C_P(z) = 1$  for all  $z$ ,  $|z| \leq 1$ , and for a renewal process,  $C_P(z) =$  squared coefficient of variation, for all  $z$ ,  $|z| \leq 1$ .  $C_P(z)$  can be considered as a index which shows how similar the considered process is to a Poisson process.

In particular, covariance between two consecutive interdeparture times  $c_1$  is,

$$c_1 = \mathbf{x}_0 R^{-1}(0) U_1(0) A_0(0) R^{-1}(0) \mathbf{e} - (\mathbf{x}_0 R^{-1}(0) \mathbf{e})^2, \quad (5.18)$$

which shows that the waiting room size affects  $c_1$  only through the idle distribution  $\mathbf{x}_0$ .

This fact is valid for  $c_k$ ,  $k < K$ .



## 6. Counting process

In this section, we consider the distribution of departure counts for  $t$  from an arbitrary departure epoch. Set the origin of the time axis at a departure epoch and let  $N(t)$  be the number of departures during  $(0, t]$ . Note [Cox 66]

$$\Pr(t_1 + \dots + t_n > t) = \Pr(N(t) < n).$$

Thus, we can obtain

$$\begin{aligned} N(z, s) &\triangleq \sum_{k=0}^{\infty} \int_0^{\infty} e^{-st} \Pr(N(t) = k) dt z^k \\ &= \frac{1}{s} + \frac{1 - z^{-1}}{s} D(z, s). \end{aligned} \quad (6.1)$$

Equation (6.1) provides the moments of the number of departures. By considering the formulas obtained by differentiating Eq. (6.1), we get

$$\int_0^{\infty} e^{-st} E[N(t)] dt = \frac{1}{s} D(1, s) \quad (6.2),$$

$$\int_0^{\infty} e^{-st} E[N^2(t)] dt = -\frac{1}{s} D(1, s) + \frac{2}{s} \left[ \frac{\partial}{\partial z} D(z, s) \right]_{z=1}. \quad (6.3)$$

From Eq. (4.8), the LSTs of  $E[N(t)]$  and  $E[N^2(t)]$  are obtained. In particular, for an N/D/1 queue, from Eq. (5.8),

$$\begin{aligned} &\int_0^{\infty} e^{-st} E[N(t)] dt \\ &= \frac{e^{-sh}}{s(1 - e^{-sh})} - \frac{e^{-sh}}{(1 - e^{-sh})^2} \mathbf{x}_0(I - L(e^{-sh})) \{sI - R(0)(I - L(e^{-sh}))\}^{-1} \mathbf{e} \quad (6.4) \\ &\int_0^{\infty} e^{-st} E[N^2(t)] dt \\ &= \frac{e^{-sh} + e^{-2sh}}{s(1 - e^{-sh})^2} - \frac{e^{-sh} + 3e^{-2sh}}{(1 - e^{-sh})^3} \mathbf{x}_0(I - L(e^{-sh})) \{sI - R(0)(I - L(e^{-sh}))\}^{-1} \mathbf{e} \end{aligned}$$

$$\begin{aligned} &+ 2 \frac{e^{-2sh}}{(1 - e^{-sh})^2} \mathbf{x}_0 \left[ \frac{\partial}{\partial z} L(z) \right]_{z=e^{-sh}} \{sI - R(0)(I - L(e^{-sh}))\}^{-1} \mathbf{e} \\ &+ 2 \frac{e^{-2sh}}{(1 - e^{-sh})^2} \mathbf{x}_0(I - L(e^{-sh})) \{sI - R(0)(I - L(e^{-sh}))\}^{-1} \\ &R(0) \left[ \frac{\partial}{\partial z} L(z) \right]_{z=e^{-sh}} \{sI - R(0)(I - L(e^{-sh}))\}^{-1} \mathbf{e}. \end{aligned} \quad (6.5)$$



## 7. Examples

### 7.1 Departure process from an M/M/1 queue

We assume that customers arrive in a Poisson process at rate  $\lambda$  and that a service time is exponentially distributed at mean  $\mu^{-1}$ , and that there are infinite waiting rooms. For an infinite system, Eqs. (3.5), (3.6) and (3.8) reduce to

$$G^{[i,i']}(z,s) = (G^{[1,0]}(z,s))^{i-i'} \quad i \geq i' \geq 0, \quad (7.1)$$

$$G^{[1,0]}(z,s) = z \sum_{k=0}^{\infty} A_k(s) (G^{[1,0]}(z,s))^k, \quad (7.2)$$

and

$$W(z,s) = (I - \sum_{k=0}^{\infty} U_{k+1}(s) G^{[k+1,0]}(z,s))^{-1}. \quad (7.3)$$

Apply the results in Section 4 with

$$\begin{aligned} m &= 1 \\ H(s) &= \frac{\mu}{s + \mu} \\ R(z) &= (z - 1)\lambda \\ S(s) &= -\frac{s}{s + \lambda} \\ x_k &= (1 - \lambda/\mu)(\lambda/\mu)^k \quad k \geq 0 \\ U_k(s) &= \begin{cases} \frac{\lambda}{s + \lambda} & k = 1 \\ 0 & k \geq 2 \end{cases} \end{aligned}$$

Thus, we obtain from Eqs. (7.2) and (7.3)

$$G^{[1,0]}(z,s) = \frac{\mu z}{s + \lambda + \mu - \lambda G^{[1,0]}(z,s)} \quad (7.4)$$

$$W(z,s) = (1 - \frac{\lambda}{s + \lambda} G^{[1,0]}(z,s))^{-1}. \quad (7.5)$$

Therefore, from Eq. (4.8),

$$D(z,s) = \frac{\lambda z}{s - \lambda z + \lambda} \quad (7.6)$$

$$= \sum_{k=1}^{\infty} \left( \frac{\lambda}{s + \lambda} \right)^k z^k. \quad (7.7)$$

Equation (7.7) means that the sum of consecutive  $k$  interdeparture times is the sum of  $k$  exponential distributed random variables. Furthermore, in view of Eq. (6.1),

$$N(z,s) = \frac{1}{s - \lambda z + \lambda} \quad (7.8)$$

$$= \int_0^{\infty} e^{-st} \sum_{k=0}^{\infty} e^{-\lambda t} \frac{(\lambda t)^k}{k!} z^k. \quad (7.9)$$

Equation (7.9) means that the distribution of the number of departures is Poisson. These results are consistent with Burke [Burke 56].

### 7.2 Departure process from an MMPP/D/1 queue

We assume that customers arrive in a Markov modulated Poisson process at a single server queue with infinite waiting rooms. For an MMPP,

$$R(z) = (z - 1)\Delta(\lambda) + Q^* \quad (7.10)$$

$$U_k(s) = \begin{cases} (sI - R(0))^{-1} \Delta(\lambda) & k = 1 \\ 0 & k \geq 2. \end{cases} \quad (7.11)$$

Applying Eq. (5.17) yields the numerical results shown in Figs. 1 and 2. In these figures,  $C_P(z)$ , the index of similarity to a Poisson process defined in Section 5, is plotted. Dotted lines show  $C_P(z)$  for an arrival process.

In Figure 1,  $m = 2$  and

$$Q^* = \begin{pmatrix} -0.1 & 0.1 \\ 0.1 & -0.1 \end{pmatrix}.$$



The shape of a curve for a departure process is similar to that for a corresponding arrival process, while  $C_P(0)$  is quite different.  $C_P(0)$  is the squared coefficient of variation and the shape of a curve is determined by the mean interdeparture time and covariances  $c_1, c_2, \dots$ . Thus, it is concluded that the covariance structure may be preserved even when the variance decreases.

In Figure 2,  $\lambda_1 = 1.0, \lambda_2 = 0.5, m = 2$  and

$$Q^* = \begin{pmatrix} -\tau & \tau \\ 0.1 & -0.1 \end{pmatrix}.$$

Here,  $C_P(z)$  is plotted as a function of  $\tau$ . It is observed again that the shape of curves of arrival and departure processes are similar.

### 7.3 The departure from a packetized voice multiplexer

In this subsection, our result is applied to the analysis of the interdeparture times of a voice packet multiplexer. Consider a multiplexer with 100 off-hook users or with 50 off-hook users. The packet arrival process from a single voice source consists of arrivals occurring at fixed intervals of 16 ms during talkspurts and no arrivals during silences. The length of a talkspurt and a silence period of a single voice source are assumed to be exponentially distributed with means 352 ms and 650 ms, respectively [Heffes 86, Sriram 86]. The packet size is assumed to be 64 bytes.

Furthermore, we adopt the successful modelling in [Heffes 86] for the superposed arrival process, which is the 2-state MMPP fitted with the following characteristics:

- 1) the mean arrival rate;
- 2) the variance-to-mean ratio of the number of arrivals in  $(0, T_1)$ ;

3) the long term variance-to-mean ratio of the number of arrivals; and

4) the third moment of the number of arrivals in  $(0, T_2)$ .

We use  $T_1 = T_2 = 0.5$ s in this paper.

Figure 3 shows  $C_P(z)$  of the arrival and departure processes. The departure process is fairly smooth for 100 users. While the Poisson approximation for the arrival process will give an optimistic result, that for the departure process will be too pessimistic and  $E_2$  can approximate the departure process if we take the result in the Appendix into account.

For 50 users, the departure process is less smooth. This is because for light load most of arriving voice packets are transmitted without waiting.

### Conclusions

The departure process of an N/G/1 queue was investigated. The first passage analysis yields LSTs of distributions of the interdeparture times.



## Appendix. $C_P(z)$ characterization

This Appendix describes a method for obtaining performance measures of as a single server queue when we can evaluate  $C_P(z)$  of its arrival process. We use the argument similar to [Sriram 88] and derive an renewal approximation based on  $C_P(z)$ .

Consider a single server queue and its arrival process with  $C_P(z)$ . Suppose the usage of the server is  $\rho$ . When an arriving customer sees  $k$  customers waiting, the  $k$  consecutive intervals of the arrival process directly interact with each other and correlations among  $k$  consecutive intervals play an important role in the queueing behavior. We take only this direct interaction into consideration, and let  $p_a(k)$  be the stationary probability that an arriving customer see  $k$  customers waiting in the system. Thus, with the probability  $p_a(k)$ , variation of  $\tau_1 + \dots + \tau_k$  contributes to the queueing behavior, where  $\tau_k$  is the  $k$ -th interarrival time. In other words, interaction between  $\tau_1$  and  $\tau_k$  influences the queueing behavior with probability  $\sum_{j=k}^{\infty} p_a(j)$ .

We approximate  $p_a(j)$  by the probability that an arriving customer see  $k$  customers waiting in an M/M/1 queue with usage  $\rho$ ,

$$p_a(j) \approx (1 - \rho)\rho^j \quad j \geq 0. \quad (A1)$$

Thus, covariance between  $\tau_1$  and  $\tau_k$ ,  $c_{k-1} = \text{cov}(\tau_1, \tau_k)$ , contributes directly to the queueing performance measures with probability  $\sum_{j=k}^{\infty} p_a(j) \approx \rho^j$ , and

$$C_P(\rho) = \sum_{j=0}^{\infty} c_k \rho^k / (\text{the mean interarrival time})^2$$

can be regarded as aggregation of such contributions. We identify the renewal process with the same  $C_P(\rho)$  as the arrival process considered. Since  $C_P(z)$  of a renewal process is

squared coefficient of variation for all  $z$ ,  $|z| < 1$ , our approximation implies that the arrival process is approximated by the renewal process whose squared coefficient of variation is equal to  $C_P(\rho)$  of the arrival process.

An example is presented for evaluating the validity of  $C_P(\rho)$  approximation. Consider tandem queues, MMPP/D/1  $\rightarrow$  ·/G/1. Assume that the mean service times of both servers are 1 and that both queues have infinite waiting rooms. For the arrival process to the first queue, we suppose  $m = 2$  and

$$Q^* = \begin{pmatrix} -0.1 & 0.1 \\ 0.1 & -0.1 \end{pmatrix}.$$

$C_P(\rho)$  of the arrival process to the second queue, that is, the departure process from the first queue was derived in 7.2 and is available. An approximation formula [Krämer 76] for a GI/G/1 queue is employed. The mean number of customers in the second queue is shown in Figure 4 and comparisons are made by simulation. Three curves are plotted for hyperexponential ( $H_2$ ), exponential ( $M$ ) and Erlang ( $E_4$ ) service in the second queue. Our approximation gives a good estimation of the mean number of customers and  $C_P(z)$  can characterize the arrival process.



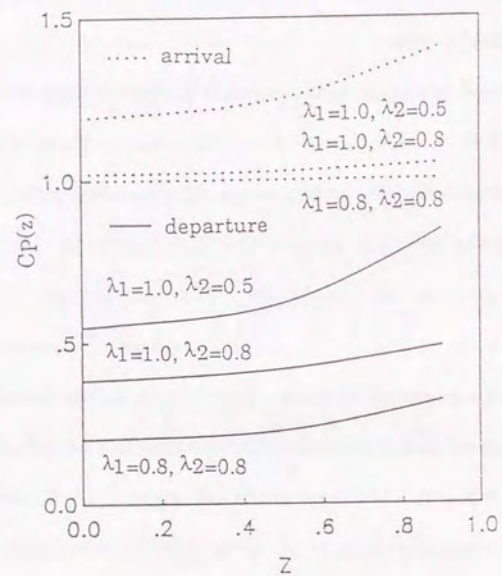


Figure 1. Departure from MMPP/D/1 queues

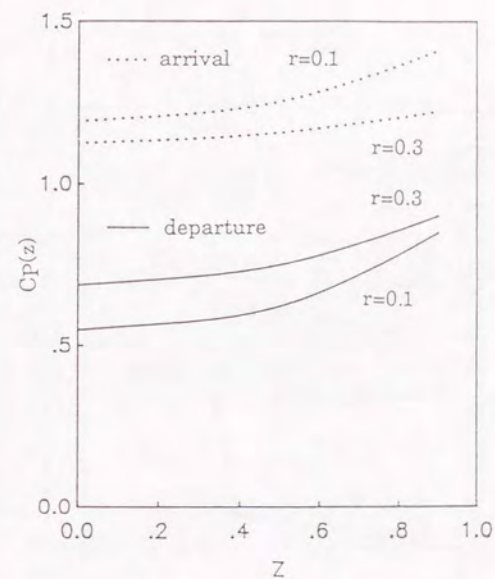


Figure 2. Departure from MMPP/D/1 queues



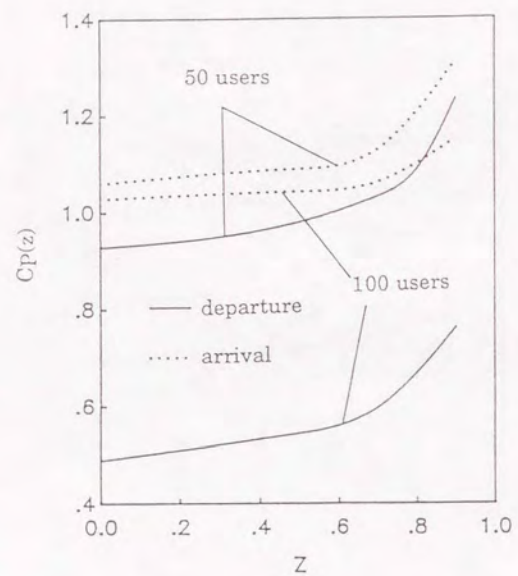


Figure 3. Departure from a packetized voice multiplexer

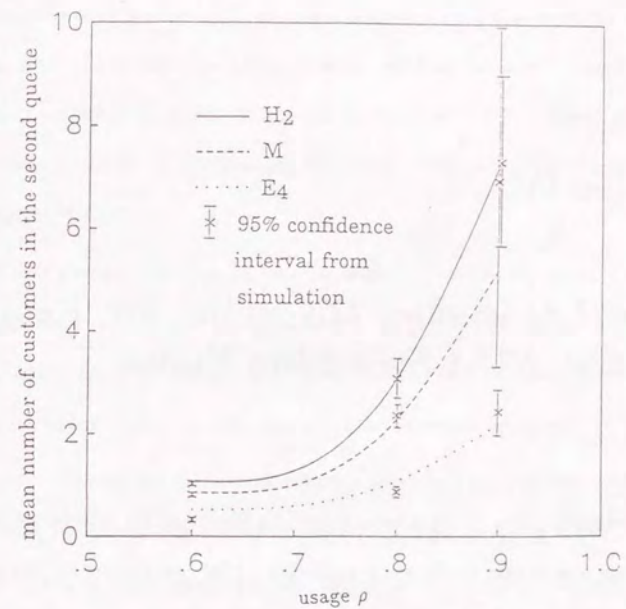


Figure 4. Mean number of customers in the second queue



## Chapter IV

### Optimal Queueing Discipline for Real-Time Traffic at ATM Switching Nodes

*Queueing disciplines at asynchronous transfer mode (ATM) switching nodes handling various kinds of real-time traffic are investigated. ATM can support various new services including voice, data and video. However, the characteristics of superposed traffic carried by ATM are not known, and a control effective for a versatile arrival process is required. This chapter derives the optimal discipline which minimizes the number of cells being delayed beyond the specified maximum allowable time and thus being discarded, without assumptions on the arrival process of cells and buffer management schemes. It also discusses implementation of the optimal discipline and a method of satisfying cell loss probability requirements of individual classes. [Saito 90a]*

#### 1. Introduction

Integration of various services such as voice, video and data into an ISDN has been of considerable interest recently. Switching systems based on the asynchronous transfer mode (ATM) are expected to provide a base for an advanced ISDN [Thomas 84, Gonet 86, Muise 86, Turner 86b, Boyer 87, Eklundh 88a,b, Saito 90c], and CCITT is working on standardizing ATM [CCITT 87].

Although this new technique seems very promising, there are still some design and performance problems to be solved [Kawashima 89]. A queueing discipline in ATM networks is one problem under discussion. There is controversy over the tradeoff between the complexity of cell handling and the effectiveness of the queueing discipline.

This chapter addresses the problem of queueing disciplines for real-time traffic. ATM switching systems have to flexibly and effectively switch many classes of traffic, including real-time traffic, such as voice, video and high-fidelity stereo. Real-time communication's grade of service greatly depends on delay. For example, the candidate values of maximum allowable end-to-end delay for voice are 100 ms [Kirton 87] or 200 ms [Gruber 83]. Cells delayed beyond these values may be discarded at the receiver or transit nodes.

The optimization problem of minimizing the number of cells discarded because of excess delay is formulated and solved.

We briefly survey here the optimal controls in integrated communications networks. It is well known that a preemptive priority discipline can minimize the weighted sum of the mean waiting time of each class in an M/M/1 queue and that head of the line priority among non-preemptive disciplines can minimize it in an M/G/1 queue [Gelenbe 80, p.199],



and later it was extended in [Buyukkoc 85, Ross 88]. Congestion control using embedded coding [Bially 80a] which maximizes the mean bit rate of voice with a constraint of voice and data packet delay, has a simple structure: a feedback with saturation [Saito 89a]. In hybrid switching systems, window flow control maximizes throughput with a constraint on the mean delay of data traffic [Vakil 87].

The optimal control derived in this chapter depend on neither the cell arrival processes, buffer management schemes [Irland 87] nor buffer size. Thus, optimal control is robust in many respects. In addition, a performance limit on the number of cells discarded because of excess delay can be estimated, if the optimal control is known. Implementation and a quality control which guarantees the specified cell loss probabilities for individual classes are also discussed.

## 2. Model

Consider a transmission link of an ATM switching node handling  $n$  classes of real-time traffic. Class  $i$  traffic permits the maximum delay  $T_i$ . Cells which exceed this maximum are assumed to be discarded.

A cell has a fixed duration and requires a unit of time for transmission. One empty or non-empty cell is always transmitted every unit of time [Devault 88]. Thus, the time axis is divided into slots one unit of time long. We assume that a cell arrives at the beginning of a slot, and leaves or is discarded at the end of a slot. The decision epoch is set just after arrival, at the beginning of a slot (Figure 1). The criterion adopted here for the queueing discipline is the number of cells discarded because of delay longer than the specified maximum.

No statistical assumption for the cell arrival process is needed. The following arguments are valid for any size of buffers, finite or infinite, and for any buffer management scheme (e.g., common or individual buffer systems [Irland 87, Köner 83]).

## 3. The optimal queueing discipline

Let  $d(t; \mathbf{u})$  be the number of cells due to be transmitted by  $t$ , but not transmitted by  $t$ , and so discarded under control  $\mathbf{u}$  at  $t$ . Let  $D(t; \mathbf{u})$  be the accumulated number of cells discarded during  $[0, t]$ . That is,

$$D(t; \mathbf{u}) = \sum_{s=0}^t d(s; \mathbf{u}). \quad (3.1)$$

Consider the optimization problem over the  $T$  horizon,

$$\min_{\mathbf{u}} D(T; \mathbf{u}). \quad (3.2)$$

For criterion (3.2), we obtain the optimal queueing discipline called the Earliest-Due-Date (EDD). The due date of a class  $i$  cell is  $T_i$  added to the arrival time. The cell with the earliest due date is transmitted if the queue discipline is EDD. Since the EDD discipline preserves the order of cells within a class, resequencing at receiving nodes is not necessary.

**Proposition.** The Earliest-Due-Date discipline is optimal for criterion (3.2).

The proof is in the Appendix.

The EDD discipline, sometimes called dynamic priority [Jaiswal 68, p. 198], was introduced by Jackson [Jackson 60,61,62], and later, analyzed in [Holtzman 71, Goldberg 77, Netterman 79, Bagchi 85].

It is known that the EDD discipline minimizes the maximum job lateness which equals service completion time minus due date, and tardiness, or  $\max(0, \text{lateness})$ , for the problem



of scheduling jobs [Conway 67, p. 30]. It is also known that EDD can satisfy any relative average delay requirements for different classes and can equalize the lateness distribution tails for different classes [Lim 88]. An implementation of the EDD discipline has been proposed [Lim 88]. Sufficient conditions for guaranteeing that all voice packets meet their due date, are discussed in [Arthurs 79]. The tree random access protocol proposed by Kurose et al. [Kurose 88] uses due date implicitly and a collision is resolved by giving transmission right to the set of stations with earlier due dates.

A similar problem has recently been solved independently in [Panwar 89] with a different approach. In [Panwar 89], the optimal discipline is called the shortest time to extinction policy.

#### 4. Implementation

The EDD discipline can be implemented in several ways. One method employing the time stamp technique [Montgomery 83] involves time-stamping the due date on the cell headers. The switch node then transmits cells in the order of stamping. Another method uses time-stamping of the maximum permissible delay at the source node and subtracting the actual delay from the time stamp at each node [Montgomery 83]. In other words, stamped time on a cell header is the slack time (the amount of time remaining before the due date) of the cell. The time-stamped cells are put in a queue in the order of stampings, regardless of cell class, and are served in an FCFS order. Thus, the number of service classes and the maximum allowable delay can be changed without changing the processing at nodes. As a result, the flexibility of the network increases and it becomes possible to fulfill the various requirements of users. The explicit relationship between the target grade

of service and the discipline seems advantageous.

Although implementation by the time stamp technique is attractive, the extra processing time required for reading and writing a time stamp on a cell is crucial. Thus, in [Lim 88], the EDD discipline is called the Head-of-the-Line with Priority Jumps discipline, and an implementation method that does not use a time stamp is developed. There are  $n$  queues in that system, and a list of cell arrival times and a clock are required for each queue.

Here, implementation using shift registers is proposed (Figure 2). When a class  $i$  cell arrives, it is put in the  $T_i$ -th slot of a shift register. A register shifts when a cell (empty or non-empty) is transmitted. Therefore, the address of a register occupied by a cell is the slack time of the cell. The cell in the nearest to the output (that is, the cell with the minimum slack time, which means the earliest due date) is transmitted.

In Figure 2, the class 1 cell in the first slot is transmitted first. Next, the cell in the third slot of class 2 is transmitted as soon as the transmission of the previous cell is completed, since the second slots of the registers are empty.



## 5. Cell loss probability of each class

In ATM systems, cell loss probability is one of the main measures of quality and is considered to be specified for each class. Cells are lost because of buffer overflow as well as by being discarded because of excess delay. While the EDD discipline can minimize the total cell loss probability on account of passing the allowable delay and transmitted cells satisfy the specific permissible delay of each class, it does not guarantee satisfaction of the specified cell loss probability for each class. It is necessary to develop a quality control which allots the total loss probability into individual classes, combined with a congestion control, or a buffer management scheme which provides a feasible region of individual cell loss probabilities, if necessary.

Under the EDD discipline, only the due date of a cell is considered, not the class to which the cell belongs. Thus, whichever cell may be transmitted, the control holds optimal, only if the transmitted cell has the earliest due date.

When there is more than one cell with the same due date (the same slack time), a quality control method in which a class  $i$  cell is transmitted first with probability  $\alpha_i$ , is expected to effectively satisfy the cell loss probability requirement of each class. Here,  $\{\alpha_i\}$  is determined appropriately to take quality requirements into account.

For example, there are two cell classes in the system. The number of arriving cells within a class is assumed to be at most one during a unit of time. If we can employ a quality control method in which a class 1 cell is transmitted first with probability  $\alpha_1$  when a class 1 cell and a class 2 cell have the same due date, we can arbitrarily set the proportion of the class one cells (or class two cells) with excess delay among whole cells with excess

delay (See 6. A numerical example). This is because a cell is discarded only when both slack times of a class 1 cell and a class 2 cell become 1 simultaneously, and a class 1 cell is discarded with probability  $1 - \alpha_1$ . Here, the number of whole classes of cells with excess delay or lost on account of buffer overflow, is determined by a congestion control, a buffer management scheme and buffer size.

When the quality control based on the EDD discipline cannot attain a quality requirement on a loss probability for each class for any tuning parameters  $\{\alpha_i\}$ , buffer size should be increased or, a typical congestion control for ATM, or a buffer management scheme [Irland 87], [Köner 83] must be introduced.

Only a cell which is not regulated under a congestion control and finds an unoccupied buffer, is subject to the quality control. The quality control method employed with a properly chosen congestion control or buffer management scheme and a proper buffer size is expected to satisfy the specified loss probabilities for individual cell classes. Even when such controls are used together, the EDD discipline for the arrival processes of cells which are accepted in the system, is optimal and minimizes the number of cells with excess delay. An example is given in the following section showing that EDD quality control is better than any other control not only in the total number of cells discarded but also in the number of individual classes of cells discarded.



## 6. A numerical example

The effectiveness of the quality control method based on the EDD discipline is verified by simulation. We assume that there are 2 cell classes. Suppose that the number of class  $i$  ( $i = 1, 2$ ) cells arriving during a slot is at most one, and let  $\lambda_i$ ,  $i = 1, 2$ , be the probability of a class  $i$  cell arriving during a slot.

If a class 1 cell and a class 2 cell have the same due date, then a class 1 cell is transmitted first with probability  $\alpha_1$ . If  $T_1 = 10, T_2 = 20$  and there are infinite buffers, simulation runs  $10^5$  units of time, and we obtain the following results. Here, 'priority' denotes the discipline in which a class 1 cell has priority over a class 2 cell and a class  $i$  cell which exceeds  $T_i$  is discarded.

Table 1 shows the effectiveness of the EDD in total cell loss probability, compared with a priority scheme. If we set  $\alpha_1 = 1.0$ , the EDD attains the loss probability of class 1 cells = 0 and that of class 2 cells =  $3.88 \times 10^{-2}$ . Thus, the EDD quality control is better than the priority scheme also in individual loss probabilities. Actually, the quality control based on the EDD discipline is superior to any control also in all of the individual loss probabilities in this example, because the total number of cells discarded is minimized and the proportion of class 1 cells discarded is arbitrarily set by  $\alpha_1$ .

As stated in the previous section, the loss probability of each class can satisfy any specified loss probability by choosing an appropriate  $\alpha_1$ , based on the assumption of this example, only if the total loss probability =  $1.55 \times 10^{-2}$ . However, if a requirement of loss probability for each class results in a requirement of total loss probability  $< 1.55 \times 10^{-2}$ , congestion control is necessary.

## 7. Conclusions

Queueing disciplines at ATM switching nodes for handling various kinds of real-time traffic were investigated. It was shown that the EDD discipline minimizes the number of cells exceeding the specified maximum allowable delay and thus discarded. Under the EDD discipline, the transmitted cells satisfy the delay quality requirement of each class. The implementation of the EDD discipline employing shift registers was considered.

Since the EDD discipline does not guarantee satisfaction of the loss probability of each class, a quality control based on the EDD discipline was proposed. The satisfaction of loss probability for each class seems attainable by combining an appropriate buffer management scheme or congestion control with the EDD. The optimality of the EDD would not be lost in that case, since the argument stated here does not depend on the buffer management scheme, buffer size or an arrival process. In addition, the optimality of the EDD for any arrival process with unknown statistics is an advantage for a queueing discipline to be employed in future ISDNs.

Performance measures for the EDD discipline in the systems including non-real-time traffic as well as real-time traffic need further study.



# Appendix: Proof of the proposition

To simplify the notation, we assume that the number of cell classes is 2. Without loss of generality, we can assume  $T_2 \geq T_1$ .

Let  $m_i^{(j)}(t; \mathbf{u})$  be the number of class  $i$  cells with slack time (= the amount of time remaining before the due date)  $j$  at  $t$  under control  $\mathbf{u}$ , where  $1 \leq j \leq T_i$ . Let  $m^{(j)}(t; \mathbf{u})$  be the number of cells with slack time  $j$  at  $t$  under control  $\mathbf{u}$ , that is,

$$\begin{cases} m^{(j)}(t; \mathbf{u}) = m_1^{(j)}(t; \mathbf{u}) + m_2^{(j)}(t; \mathbf{u}) & 1 \leq j \leq T_1, \\ m^{(j)}(t; \mathbf{u}) = m_2^{(j)}(t; \mathbf{u}) & T_1 < j \leq T_2. \end{cases} \quad (A.1)$$

$\{z_k(t; \mathbf{u}), 0 \leq k \leq T_2\}$  is the number of cells with slack time less than  $k+1$  at  $t$  under control  $\mathbf{u}$  added to the number of discarded cells during  $[0, t]$  under control  $\mathbf{u}$ .

$$\begin{cases} z_0(t; \mathbf{u}) = D(t; \mathbf{u}) \\ z_k(t; \mathbf{u}) = D(t; \mathbf{u}) + \sum_{j=1}^k m^{(j)}(t; \mathbf{u}) & 1 \leq k \leq T_2 \end{cases} \quad (A.2)$$

Here,  $D(t; \mathbf{u})$  is the number of discarded cells during  $[0, t]$  and is defined in Section 3.

Consider the evolution of  $\{z_k(t; \mathbf{u}), 0 \leq k \leq T_2\}$ . Let  $S_{\mathbf{u}}$  be the slack time of the cell transmitted at  $t$  under control  $\mathbf{u}$ . (For notational convenience, denote  $S_{\mathbf{u}} = T_2 + 1$  that an empty cell is transmitted at  $t$  under control  $\mathbf{u}$ .) Cells which are not transmitted with slack time 1 are discarded, and the slack time of cells which are not transmitted with slack time  $j+1$  at  $t$  are  $j$  at  $t+1$ . Therefore,

$$\begin{aligned} D(t+1; \mathbf{u}) &= D(t; \mathbf{u}) + m^{(1)}(t; \mathbf{u}) - \mathbf{1}(S_{\mathbf{u}} = 1) \\ m^{(j)}(t+1; \mathbf{u}) &= \begin{cases} m^{(j+1)}(t; \mathbf{u}) - \mathbf{1}(S_{\mathbf{u}} = j+1) & \text{for } j \neq T_1, T_2 \\ m^{(j+1)}(t; \mathbf{u}) + a_1(t+1) - \mathbf{1}(S_{\mathbf{u}} = j+1) & \text{for } j = T_1 \\ a_2(t+1) & \text{for } j = T_2 \end{cases} \end{aligned}$$

Here,  $a_i(t+1)$ ,  $i = 1, 2$ , is the number of class  $i$  cells arriving at  $t+1$ , and  $\mathbf{1}(\cdot)$  is an indicator function,

$$\mathbf{1}(z) = \begin{cases} 1, & \text{if } z \text{ is true} \\ 0, & \text{if } z \text{ is false.} \end{cases}$$

Thus, we obtain,

$$z_0(t+1; \mathbf{u}) = z_1(t; \mathbf{u}) - \mathbf{1}(S_{\mathbf{u}} = 1) \quad (A.3)$$

$$\begin{aligned} z_k(t+1; \mathbf{u}) &= D(t+1; \mathbf{u}) + \sum_{j=1}^k m^{(j)}(t+1; \mathbf{u}) \\ &= \begin{cases} D(t; \mathbf{u}) + m^{(1)}(t; \mathbf{u}) + \sum_{j=1}^k m^{(j+1)}(t; \mathbf{u}) - \mathbf{1}(S_{\mathbf{u}} \leq k+1) & 1 \leq k < T_1 \\ D(t; \mathbf{u}) + m^{(1)}(t; \mathbf{u}) + \sum_{j=1}^k m^{(j+1)}(t; \mathbf{u}) + a_1(t+1) - \mathbf{1}(S_{\mathbf{u}} \leq k+1) & T_1 \leq k < T_2 \\ D(t; \mathbf{u}) + m^{(1)}(t; \mathbf{u}) + \sum_{j=1}^{T_2-1} m^{(j+1)}(t; \mathbf{u}) + a_1(t+1) + a_2(t+1) - \mathbf{1}(S_{\mathbf{u}} \leq T_2) & k = T_2 \end{cases} \\ &= \begin{cases} z_{k+1}(t; \mathbf{u}) - \mathbf{1}(S_{\mathbf{u}} \leq k+1) & 1 \leq k < T_1 \\ z_{k+1}(t; \mathbf{u}) + a_1(t+1) - \mathbf{1}(S_{\mathbf{u}} \leq k+1) & T_1 \leq k < T_2 \\ z_{T_2}(t; \mathbf{u}) + a_1(t+1) + a_2(t+1) - \mathbf{1}(S_{\mathbf{u}} \leq T_2) & k = T_2 \end{cases} \quad (A.4) \end{aligned}$$

Let  $\mathbf{u}^*$  denote the EDD queueing discipline and  $S_{\mathbf{u}}^*$  be the slack time of the cell transmitted at  $t$  under the EDD discipline. Set  $S_{\mathbf{u}}^* = T_2 + 1$ , if an empty cell is transmitted at  $t$  under the EDD discipline.

For  $0 \leq k < T_2$ ,

$$\begin{aligned} &z_k(t+1; \mathbf{u}) - z_k(t+1; \mathbf{u}^*) \\ &= z_{k+1}(t; \mathbf{u}) - z_{k+1}(t; \mathbf{u}^*) - \mathbf{1}(S_{\mathbf{u}} \leq k+1) + \mathbf{1}(S_{\mathbf{u}}^* \leq k+1), \end{aligned} \quad (A.5)$$

and

$$\begin{aligned} &z_{T_2}(t+1; \mathbf{u}) - z_{T_2}(t+1; \mathbf{u}^*) \\ &= z_{T_2}(t; \mathbf{u}) - z_{T_2}(t; \mathbf{u}^*) - \mathbf{1}(S_{\mathbf{u}} \leq T_2) + \mathbf{1}(S_{\mathbf{u}}^* \leq T_2). \end{aligned} \quad (A.6)$$

We first derive the following Lemmas as preliminaries for the proof of the main theorem.



**Lemma 1.**

For any integer  $k$ ,  $0 \leq k < S_u^*$ ,

$$z_k(t; u^*) = z_0(t; u^*). \quad (A.7)$$

*Proof:* There is no cell with slack time less than  $S_u^*$  at  $t$  under  $u^*$  in the system, because a cell with the earliest due date has the minimum slack time. Therefore,

$$m^{(j)}(t; u^*) = 0, \quad 1 \leq j < S_u^*.$$

Consequently, for  $1 \leq k < S_u^*$ ,

$$\begin{aligned} z_k(t; u^*) &= D(t; u^*) + \sum_{j=1}^k m^{(j)}(t; u^*) \\ &= D(t; u^*) \\ &= z_0(t; u^*). \end{aligned}$$

For  $k = 0$ , Eq. (A.7) is valid. This completes the proof. ■

**Lemma 2.**

If  $S_u \leq T_2$ ,

$$z_j(t; u) \geq z_0(t; u) + 1 \quad \text{for } j \geq S_u. \quad (A.8)$$

*Proof:* It is trivial that

$$m^{(S_u)}(t; u) \geq 1$$

and that for any integers  $k$  and  $l$ ,  $k \geq l$ ,

$$z_k(t; u) \geq z_l(t; u).$$

Therefore, for any integer  $j \geq S_u$ ,

$$\begin{aligned} z_j(t; u) &\geq z_{S_u}(t; u) \\ &= z_{S_u-1}(t; u) + m^{(S_u)}(t; u) \\ &\geq z_0(t; u) + m^{(S_u)}(t; u) \\ &\geq z_0(t; u) + 1. \quad \blacksquare \end{aligned}$$

**Lemma 3.**

If  $S_u^* > S_u$  and  $S_u^* \leq T_2$ ,

$$\begin{aligned} &z_{T_2}(t+1; u) - z_{T_2}(t+1; u^*) \\ &= z_{T_2}(t; u) - z_{T_2}(t; u^*). \end{aligned} \quad (A.9)$$

If  $S_u^* > S_u$  and  $S_u^* > T_2$ ,

$$\begin{aligned} &z_{T_2}(t+1; u) - z_{T_2}(t+1; u^*) \\ &\geq z_0(t; u) - z_0(t; u^*). \end{aligned} \quad (A.10)$$

*Proof:* Suppose  $S_u < S_u^*$ . Since  $S_u < S_u^* \leq T_2 + 1$ ,  $S_u \leq T_2$ .

Therefore, if  $S_u^* \leq T_2$ , then  $1(S_u \leq T_2) = 1(S_u^* \leq T_2)$ , and if  $S_u^* > T_2$ , then  $1(S_u \leq T_2) - 1(S_u^* \leq T_2) = 1$ . Hence, employ Eq. (A.6),

$$\begin{aligned} &z_{T_2}(t+1; u) - z_{T_2}(t+1; u^*) \\ &= \begin{cases} z_{T_2}(t; u) - z_{T_2}(t; u^*) & \text{for } S_u^* \leq T_2 \\ z_{T_2}(t; u) - z_{T_2}(t; u^*) - 1 & \text{for } S_u^* > T_2. \end{cases} \end{aligned} \quad (A.11)$$

The proof of Eq. (A.9) is completed.

If  $S_u^* > T_2$ , from Lemma 1,

$$z_{T_2}(t; u^*) = z_0(t; u^*).$$



Since  $S_u \leq T_2$ , from Lemma 2,

$$z_{T_2}(t; u) \geq z_0(t; u) + 1.$$

Using Eq. (A.11), for  $S_u^* > T_2$ ,

$$\begin{aligned} & z_{T_2}(t+1; u) - z_{T_2}(t+1; u^*) \\ & \geq z_0(t; u) - z_0(t; u^*). \quad \blacksquare \end{aligned}$$

**Lemma 4.**

If  $S_u^* > S_u$ , then for any integer  $k$  such that  $S_u^* \leq k+1 \leq T_2$  or  $S_u > k+1 \geq 1$ ,

$$\begin{aligned} & z_k(t+1; u) - z_k(t+1; u^*) \\ & = z_{k+1}(t; u) - z_{k+1}(t; u^*). \end{aligned}$$

*Proof:* If  $S_u^* > S_u$ , then for any integer  $k$  such that  $S_u^* \leq k+1 \leq T_2$  or  $S_u > k+1 \geq 1$ ,

$$\begin{aligned} 1(S_u \leq k+1) &= 1(S_u^* \leq k+1) \\ &= \begin{cases} 1 & \text{for } S_u^* \leq k+1 \\ 0 & \text{for } S_u^* > k+1, \end{cases} \end{aligned}$$

since  $S_u^* > S_u$ . Note that if  $S_u > k+1$ , then  $T_2 > k$ , since  $T_2 + 1 \geq S_u^* > S_u > k+1$ .

Therefore, from Eq. (A.5),

$$\begin{aligned} & z_k(t+1; u) - z_k(t+1; u^*) \\ & = z_{k+1}(t; u) - z_{k+1}(t; u^*). \quad \blacksquare \end{aligned}$$

**Lemma 5.**

If  $S_u^* > S_u$ , then for any integer  $k$  such that  $S_u \leq k+1 < S_u^*$ ,

$$\begin{aligned} & z_k(t+1; u) - z_k(t+1; u^*) \\ & \geq z_0(t; u) - z_0(t; u^*). \end{aligned}$$

*Proof:* If  $S_u^* > S_u$ , then for any integer  $k$  such that  $S_u \leq k+1 < S_u^*$ , from Eq. (A.5),

$$\begin{aligned} & z_k(t+1; u) - z_k(t+1; u^*) \\ & = z_{k+1}(t; u) - z_{k+1}(t; u^*) - 1. \end{aligned}$$

From Lemma 1, for  $0 \leq j < S_u^*$ ,  $z_j(t; u^*) = z_0(t; u^*)$ . The fact  $S_u^* > S_u$  implies  $S_u \leq T_2$ .

Hence, from Lemma 2,  $z_j(t; u) \geq z_0(t; u) + 1$  for  $j \geq S_u$ . Set  $j = k+1$ . Thus, for

$S_u^* > k+1 \geq S_u$ ,

$$\begin{aligned} & z_k(t+1; u) - z_k(t+1; u^*) \\ & \geq z_0(t; u) - z_0(t; u^*). \quad \blacksquare \end{aligned}$$

We can now state the main theorem.

**Theorem.**

For any integers  $t$  and  $k$ ,  $0 \leq t \leq T$ ,  $0 \leq k \leq T_2$ ,

$$z_k(t; u^*) \leq z_k(t; u). \quad (A.12)$$

*Proof:* The proof is by the mathematical induction on  $t$ .

For  $t = 0$ ,  $D(0; u) = D(0; u^*) = 0$  and  $m^{(j)}(0; u) = m^{(j)}(0; u^*)$  for  $1 \leq j \leq T_2$ , thus

for  $0 \leq k \leq T_2$ ,

$$z_k(0; u) = z_k(0; u^*). \quad (A.13)$$

Suppose

$$z_k(t; u) \geq z_k(t; u^*), \quad 0 \leq k \leq T_2 \quad (A.14)$$

at some  $t$ ,  $0 \leq t < T$ .

Consider the two disjoint cases, (i)  $S_u^* \leq S_u$  and (ii)  $S_u^* > S_u$ . We first complete the proof for the case (i) and then do it for the case (ii).



(i) Assume  $S_u^* \leq S_u$ . Thus, using Eqs. (A.5) and (A.6),

$$\begin{aligned} & z_k(t+1; u) - z_k(t+1; u^*) \\ & \geq \begin{cases} z_{k+1}(t; u) - z_{k+1}(t; u^*) & \text{for } 0 \leq k < T_2, \\ z_{T_2}(t; u) - z_{T_2}(t; u^*) & \text{for } k = T_2. \end{cases} \end{aligned}$$

From the assumption (A.14),

$$z_k(t+1; u) - z_k(t+1; u^*) \geq 0.$$

(ii) Assume  $S_u^* > S_u$ .

First consider the case  $k = T_2$ . From Lemma 3, for  $S_u^* \leq T_2$ ,

$$z_{T_2}(t+1; u) - z_{T_2}(t+1; u^*) = z_{T_2}(t; u) - z_{T_2}(t; u^*),$$

and for  $S_u^* > T_2$ ,

$$z_{T_2}(t+1; u) - z_{T_2}(t+1; u^*) \geq z_0(t; u) - z_0(t; u^*).$$

From assumption (A.14),

$$z_{T_2}(t+1; u) - z_{T_2}(t+1; u^*) \geq 0.$$

This completes the proof for the case  $k = T_2$ .

Next, consider the case  $0 \leq k < T_2$ . From Lemmas 4 and 5, for any integer  $k$ ,  $0 \leq k < T_2$ ,

$$\begin{aligned} & z_k(t+1; u) - z_k(t+1; u^*) \\ & = z_{k+1}(t; u) - z_{k+1}(t; u^*), \end{aligned}$$

or

$$\begin{aligned} & z_k(t+1; u) - z_k(t+1; u^*) \\ & \geq z_0(t; u) - z_0(t; u^*). \end{aligned}$$

From assumption (A.14),

$$z_k(t+1; u) - z_k(t+1; u^*) \geq 0.$$

The induction is completed. ■

The proof of the proposition is directly derived from the theorem proved above. Apply the theorem with  $t = T$  and  $k = 0$ .

$$z_0(T; u^*) \leq z_0(T; u)$$

That is, for any control  $u$ ,

$$D(T; u^*) \leq D(T; u).$$



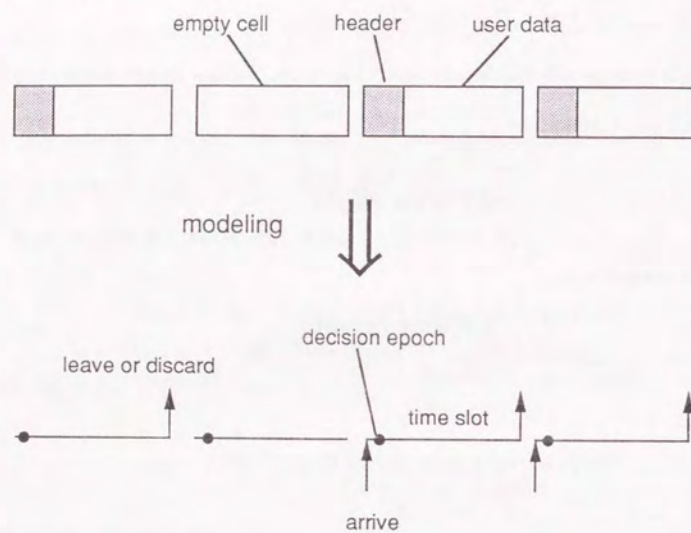


Figure 1. Modeling

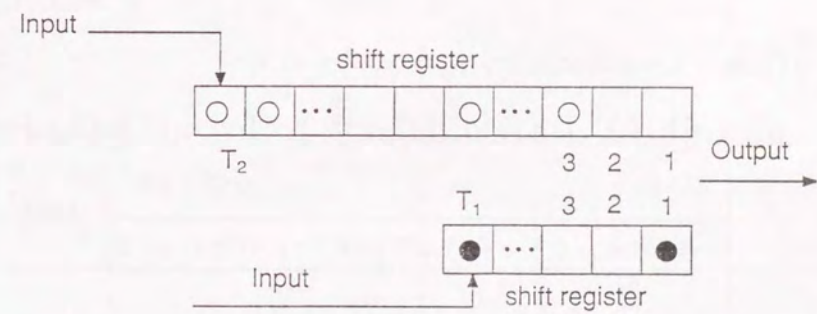


Figure 2. Implementation



Table 1. Loss Probability ( $\lambda_1=0.6, \lambda_2=0.4$ )

	class 1			class 2			total
	$\alpha_1=0.0$	$\alpha_1=0.5$	$\alpha_1=1.0$	$\alpha_1=0.0$	$\alpha_1=0.5$	$\alpha_1=1.0$	
EDD	2.57	1.28	0	0	1.94	3.88	1.55
priority	0			6.57			2.62

$\times 10^{-2}$

## Chapter V

### Optimal Control of Variable Rate Coding in ATM Networks

*The problem of controlling dynamic voice coding rate to reduce traffic congestion in an ATM network is investigated. The embedded coding scheme is considered and the design of a control is formulated as an optimization problem with a constraint. It is shown that the control which achieves the maximum average voice coding rate under an average queue length (or an average waiting time) constraint has a simple structure specified by two parameters and is a randomized modification of two feedback controls with input saturation. This result can be extended to the case that embedded coded video traffic is added to voice and data traffic. [Saito 89a]*



## 1. Introduction

Recently, there has been considerable interest in integrated networks. In particular, the development of ATM networks has become widely recognized as a significant step towards achieving an ISDN [Turner 85].

In ATM networks, a new congestion control scheme for voice cells needs to be implemented. One such scheme that has recently received attention is decreasing the bit rate of voice coding during overload [Bially 80a, Seguel 82, Listanti 83, Gafni 84, Holtzman 85, Forst 86]. Traditional voice flow control mechanisms either block the initiation of a call or discard voice cells already in progress. By contrast, variable rate coding dynamically trades off between voice quality and congestion by reducing the voice coding bit rate at the point of congestion or the point of entry. Therefore, the design objective of this control scheme is to maximize the voice quality while minimizing congestion. This problem is formulated as a voice quality optimization problem with a queue length (or a waiting time) constraint. Voice quality is known to deteriorate as coding rate decreases [Heggstad 82]. However, the effect of bit rate on voice coding is quite complicated [Cox 80, Holtzman 85], so the average voice coding rate is used as the criterion for voice quality [Bially 80a, Muise 86]. The signal-to-quantizing noise ratio with  $N$  bit quantization in many coding schemes, e.g.  $\mu$ -PCM coding, has the form,  $c_1 + c_2 N$  [Rabiner 78]. Here,  $c_1$  and  $c_2$  are constant. Therefore, the maximization of average coding rate means maximization of the segment, or cell-wise, signal-to-noise ratio.

Variable rate voice coding can be realized in a variety of ways [Bially 80a]. The first method is time stripping. This method is used for an LPC coded communication

system, which typically transmits voice in blocks every 20 ms. By stripping blocks and using suitable interpolation methods at the receiver, intelligible speech can be obtained. Another method is to strip the less significant bits from PCM (ADPCM) speech at the multiplexers or vocoders at the network entering point and at the packet switching nodes in the networks. We consider this method, which is an embedded coding scheme [Bially 80a, Goodman 80], because it has been implemented in an experimental system [Muise 86] and can be regarded as one of the most practical variable coding methods (Figure 1). Embedded coding was recommended in CCITT Recommendation G.722 in 1984 [CCITT 84]. It realizes low coding rates by reducing cell lengths [Bially 80a].

A typical example is as follows: Voice signals from off-hook users are sampled at regular intervals, typically 125  $\mu$ sec. Each sample is quantized by, for example, 4 bits. During silence periods of talkers, samples are discarded. A specified number of samples during talkspurt are grouped together, packetized and transmitted. When cells are constructed, the  $i$ -th bit of all samples are grouped to make up a cell (packet interleave method). Thus, four cells are constructed simultaneously, when a sample is quantized by 4 bits. When a network is congested, the bits of signals are dropped by discarding cells at the vocoders or multiplexers at the network entry point and at ATM nodes within the network. The cell consisting of the least significant bits is discarded first. Cells are discarded in ascending order of significance.

A considerable amount of work has been done on variable bit rate control schemes and evaluation of their performance [Bially 80a, Cox 80, Goodman 80, Seguel 82, Tham 83, Holtzman 85, Forst 86, Fredericks 86]. Gafni and Bertsekas [Gafni 84] have developed



an algorithm to allocate fairly the link capacity for regulated voice traffic.

There are several reports on applying optimization with a constraint to the design of other flow control schemes in communication networks. Maglaris and Schwartz [Maglaris 82] investigated the problem of dynamically allocating the bandwidth of a trunk to line- and packet-switched traffic. In this study, the optimality criterion is to minimize the average packet delay with the blocking probability of arriving line-switched calls constrained to no more than a specified acceptable level.

Lazar [Lazar 83a,b] and Vakil and Lazar [Vakil 87] considered optimal flow control for computer networks, and they modeled computer communication protocols with acknowledgement as closed queueing networks. They showed that the control that achieves maximum throughput under an average queue length constraint is a window flow control.

Recently, Nain and Ross [Nain 86a,b] investigated the problem of a multi-queue system competing for a single server. The optimization criterion in this problem was to minimize the linear combination of the average queue lengths with an average queue length of a queue constrained to no more than a specified level. This model can be regarded as an optimality assignment of a transmission channel for heterogeneous traffic, for example, voice and data.

This chapter proves through a Lagrange multiplier technique [Tijms 86] that the optimal control scheme has a simple structure. It is shown simultaneously that an optimal control randomizes the two deterministic controls [Nain 86a,b].

The intuitive control scheme for reducing bit rates while a queue length or a workload is beyond a specific level, has been analyzed, evaluated and implemented [Bially 80a, Seguel 82, Tham 83]. The results in the following sections show that this control scheme is not

optimal with respect to the criterion and the constraint considered in this chapter, and needs a slight modification.

## 2. Problem Formulation

Consider an ATM node in which cells arrive at the beginning of a slot and leave at the end of a slot.  $N_V$  voice cells are assumed to be generated and arrive simultaneously. We also assume that data cells may arrive in a group. Let  $N_D$  be the group size of data cells. A cell is assumed to have a fixed duration equal to the slot size. At most, one cell is transmitted in a slot. Let  $m_V$  and  $m_D$  be the number of groups of voice cells and data cells arriving at the beginning of a slot. They are assumed to be mutually independent. Define  $\lambda_V(k) = \Pr\{m_V = k\}$ ,  $\lambda_D(k) = \Pr\{m_D = k\}$ . Let  $\lambda_V$  be the average,  $m_{V2}$ , the second moment and  $M_V(z)$ , the generating function of  $m_V$ . Let  $\lambda_D$  be the average,  $m_{D2}$ , the second moment and  $M_D(z)$ , the generating function of  $m_D$ . Then, define  $h(i) = \Pr\{N_D = i\}$ , ( $i = 1, 2, \dots$ ),  $h = E[N_D]$  and  $H(z) = \sum_{i=1}^{\infty} h(i)z^i$ . Let  $\xi_D(a_D = k)$  be the probability that  $a_D$ , the total number of data cells arriving at the beginning of a slot is  $k$ , which is completely determined by  $\{\lambda_D(i), h(j)\}$ .  $N_V$  is assumed to be constant in accordance with actual systems. A group of voice cells consisting of  $N_V$  cells is reduced on arrival by discarding less significant cells [Bially 80a] according to congestion. As a result,  $N$  voice cells are accepted in a switching system. Here,  $N$  is selected as a control parameter. It is assumed that  $N_L \leq N \leq N_V$ . In other words,  $\underline{a}(i, j) = iN_L + j$  and  $\bar{a}(i, j) = iN_V + j$ , where  $\underline{a}(i, j)$  is the minimum number of accepted voice and data cells in a slot when  $m_V = i$  and  $a_D = j$ , and  $\bar{a}(i, j)$  is the maximum number of accepted voice and data cells in a slot when  $m_V = i$  and  $a_D = j$ . A non-preemptive



priority is assumed for voice cells over data cells. Let  $h$  be the average and  $h_2$  the second moment of the number of cells in a data packet (Figure 2). Notation is listed in Appendix 4.

The objective is to find a cell-dropping (bit-dropping) control  $u$  that maximizes the long-run average coding rate under an average queue length of cells constraint. The control  $u$  is restricted to the class of randomized stationary controls concerning the total number of cells in the system and the number of arriving packets in the current slot. Define  $\rho_L = \lambda_D h + \lambda_V N_L$ . If  $\rho_L \geq 1$ , there does not exist a control that meets the constraint. It is assumed that  $\rho_L < 1$  in the following, and our attention is limited to the class of controls such that there exists a stationary state under the control. Thus, the problem is formulated as

$$\underset{u}{\text{maximize}} \quad J(u) = E[N], \quad (2.1)$$

$$\underset{u}{\text{subj. to}} \quad C(u) = \sum_{k=1}^{\infty} (k-1)p_k \leq Q. \quad (2.2)$$

Here,  $p_k$  is the probability that  $x$ , the total number of cells in the system, is  $k$ . ( $E[N]$  and  $\{p_k\}$  both depend on the control policy  $u$ , but for sake of simplicity the index  $u$  has been omitted.)

Later, an average waiting time constraint instead of queue length constraint (2.2) is also considered.

### 3. Structure of the optimal control

**Lemma 1.** Our problem is equivalent to

$$\underset{u}{\text{maximize}} \quad \tilde{J}(u) = \sum_{k=1}^{\infty} q_k, \quad (3.1)$$

$$\underset{u}{\text{subj. to}} \quad \tilde{C}(u) = \sum_{k=1}^{\infty} (k-1-Q)q_k \leq Q, \quad (3.2)$$

where  $q_k = p_k/p_0$ .

*Proof:* Note

$$\lambda_V E[N] + \lambda_D h = 1 - p_0.$$

Here,  $\lambda_V$ ,  $\lambda_D$  and  $h$  do not depend on  $u$ . Therefore, maximizing  $E[N]$  is identical to minimizing  $p_0$ . Thus, using the fact that

$$p_0 = (1 + \sum_{k=1}^{\infty} q_k)^{-1},$$

(2.1) and (2.2) reduce to (3.1) and (3.2). Hence, our problem is equivalent to the optimization problem (3.1) with constraint (3.2). ■

Our focus is on the equivalent problem mentioned above. Using the Lagrange multiplier method [Nain 86a,b, Tijms 86, Ma 86], it is shown that if there is at least one control that meets the constraint, then there is an optimal control randomizing two feedback controls with input saturation.

A constrained optimization problem can be reduced to one without a constraint through the introduction of Lagrange multipliers. For each fixed multiplier  $\gamma \geq 0$ , define the Lagrangian

$$\underset{u}{\text{maximize}} \quad J_{\gamma}(u) = \sum_{k=1}^{\infty} (1 - \gamma(k-1-Q))q_k. \quad (3.3)$$



We consider the following optimization problem to be equivalent to (3.3) to simplify the notation.

$$\maximize_u J_c(u) = \sum_{k=1}^{\infty} (c-k)q_k. \quad (3.4)$$

Here,  $c = \gamma^{-1} + 1 + Q$ .

Let  $\theta_{i,j}(k,l;u)$  denote the probability under control  $u$  that  $a_{VD}$ , the total number of accepted voice and data cells in the time slot, is  $j$ , given that  $z$ , the total number of cells in the system in the preceding time slot, is  $i$ ,  $m_V$ , the number of groups of voice cells arriving at the beginning of the time slot, is  $k$  and  $a_D$ , the total number of data cells arriving at the beginning of the time slot, is  $l$ . Since attention is restricted to stationary control, control  $u$  can be specified completely by  $\{\theta_{i,j}(k,l;u), i \geq 0, j \geq 0, k \geq 0, l \geq 0\}$ . In other words, control  $u$  is the set  $\{\theta_{i,j}(k,l;u), i \geq 0, j \geq 0, k \geq 0, l \geq 0\}$ .

Here, we define the deterministic control  $u_n = \{\theta_{i,j}(k,l;u_n)\}$ , for each  $n$  ( $n = 0, 1, 2, \dots$ ). The deterministic control  $u_n$  operates such that the total number of cells reaches  $n$  if possible after accepting arriving voice and data cells. In other words,  $u_n$  is the feedback control for the difference between  $n$  and the number of cells in the system at the arrival epoch ( $= i - 1(i > 0)$ ), and saturates under the condition  $\underline{a}(k,i) \leq a_{VD} \leq \bar{a}(k,l)$ . For each fixed  $n$ , the number of accepted voice and data cells  $a_{VD}$  under control  $u_n$ , when  $z = i$ ,  $m_V = k$  and  $a_D = l$ , is

$$a_{VD} = \begin{cases} [n-i+1]_{\underline{a}(k,l)}^{\bar{a}(k,l)}, & \text{if } i > 0, \\ [n]_{\underline{a}(k,l)}^{\bar{a}(k,l)}, & \text{if } i = 0. \end{cases}$$

Here,  $y = [x]_a^b$  is a saturation function,

$$y = \begin{cases} a, & \text{if } z < a, \\ z, & \text{if } a \leq z \leq b, \\ b, & \text{if } b < z. \end{cases}$$

That is,  $u_n$  is,

$$\theta_{i,j}(k,l;u_n) = \begin{cases} 1(j = [n-i+1]_{\underline{a}(k,l)}^{\bar{a}(k,l)}), & \text{if } i > 0, \\ 1([n]_{\underline{a}(k,l)}^{\bar{a}(k,l)}), & \text{if } i = 0. \end{cases} \quad (3.5)$$

Here,  $1(\cdot)$  is an indicator function.

In the rest of this section, it is shown that the optimal control is randomization of two deterministic controls  $u_n$  and  $u_{n+1}$ .

The following lemma provides a clue to the optimal control.

**Lemma 2.1.** Consider two controls  $u = \{\theta_{i,j}(k,l;u)\}$  and  $u' = \{\theta_{i,j}(k,l;u')\}$ , where

$\theta_{i,j}(k,l;u) = \theta_{i,j}(k,l;u')$  except for  $\{(i,k,l) = (i_0,k_0,l_0)\}$ . Suppose

$$\begin{cases} \theta_{i_0,j}(k_0,l_0;u) = \theta_{i_0,j}(k_0,l_0;u'), & \text{for } j < j_0, \\ \theta_{i_0,j}(k_0,l_0;u) \leq \theta_{i_0,j}(k_0,l_0;u'), & \text{for } j_0 = j, \\ \theta_{i_0,j}(k_0,l_0;u) \geq \theta_{i_0,j}(k_0,l_0;u') = 0 & \text{for } j_0 < j. \end{cases}$$

Then,

$$\begin{cases} q_i(u) \geq q_i(u'), & \text{for } i > i_0 + j_0 - 1(i_0 \neq 0), \\ q_i(u) = q_i(u'), & \text{for } i = 0, 1, \dots, i_0 + j_0 - 1(i_0 \neq 0). \end{cases}$$

*Proof:* Consider the process of the total number of cells in the system. This process is a Markov process and the state probability  $p_i$  satisfies the equilibrium equation,

$$(1-\lambda)p_1 = \lambda p_0,$$

$$(1-\lambda)p_{i+1} = \eta_{i,2}p_i + \dots + \eta_{i,i+1}p_1 + \eta_{0,i+1}p_0 \quad i \geq 1. \quad (3.6)$$

Here,  $\lambda \triangleq 1 - \lambda_V(0)\lambda_D(0)$  and  $\eta_{i,j}$  is the probability that the total number of accepted cells is equal to or more than  $j$  when  $z = i$ ,

$$\eta_{i,j} = \sum_{m=j}^{\infty} \sum_{k=0}^{\infty} \sum_{l=0}^{\infty} \theta_{i,m}(k,l;u) \lambda_V(k) \xi_D(l).$$

Since  $q_i = p_i/p_0$ ,  $\{q_i\}$  satisfies (3.6). For controls  $u$  and  $u'$ ,

$$\begin{cases} \eta_{i,j}(u) = \eta_{i,j}(u'), & \text{for } i \neq i_0 \text{ or for } i = i_0, j \leq j_0, \\ \eta_{i,j}(u) \geq \eta_{i,j}(u'), & \text{for } i = i_0, j > j_0. \end{cases}$$



Employing (3.6),  $q_{i+1}$  can be obtained by  $q_0 (= 1), q_1, \dots, q_i$  and  $\eta_{i,2}, \dots, \eta_{i,i+1}, \eta_{0,i+1}$ .

Therefore,

$$\begin{cases} q_{i+1}(u) = q_{i+1}(u'), & \text{for } i_0 + j_0 - 1(i_0 \neq 0) \geq i + 1, \\ q_{i+1}(u) \geq q_{i+1}(u'), & \text{for } i + 1 > i_0 + j_0 - 1(i_0 \neq 0). \end{cases}$$

■

In the same way, we also obtain Lemma 2.2.

**Lemma 2.2.** Consider two controls  $u = \{\theta_{i,j}(k, l; u)\}$  and  $u' = \{\theta_{i,j}(k, l; u')\}$ , where

$\theta_{i,j}(k, l; u) = \theta_{i,j}(k, l; u')$  except for  $\{(i, k, l) = (i_0, k_0, l_0)\}$ . Suppose

$$\begin{cases} \theta_{i_0,j}(k_0, l_0; u) \geq \theta_{i_0,j}(k_0, l_0; u') = 0, & \text{for } j < j_0, \\ \theta_{i_0,j}(k_0, l_0; u) < \theta_{i_0,j}(k_0, l_0; u') = 1, & \text{for } j_0 = j, \\ \theta_{i_0,j}(k_0, l_0; u) = \theta_{i_0,j}(k_0, l_0; u') = 0, & \text{for } j_0 < j. \end{cases}$$

Then,

$$q_i(u) \leq q_i(u').$$

Proof is omitted here.

Denote  $y(i, j) \triangleq i + j - 1(i \neq 0)$ , the total number of cells in the system after accepting  $j$  voice and data cells when the total number of cells in the preceding slot =  $i$ . For each fixed  $m \geq 0$ , define  $U(m)$ , the class of the randomized stationary controls such that for all  $u = \{\theta_{i,j}(k, l; u)\} \in U(m)$ ,

$$\theta_{i,j}(k, l; u) = \begin{cases} 1, & \text{for } j = \underline{a}(k, l), \text{ where } y(i, \underline{a}(k, l)) \geq m, \\ 0, & \text{for } j \neq \underline{a}(k, l), \text{ where } y(i, \underline{a}(k, l)) \geq m, \\ 0, & \text{for } y(i, j) \geq m, \text{ where } y(i, \underline{a}(k, l)) < m \leq y(i, \bar{a}(k, l)). \end{cases}$$

$U(m)$  is the class of controls that accept minimum number of cells, if the total number of cells after accepting minimum number of cells exceeds the threshold  $m$ . Otherwise, the controls in  $U(m)$  maintain the total number of cells below  $m$ .

**Lemma 3.** Assume that there exists an integer  $K$  such that  $K - 1 < c \leq K$ . Then, the optimal control  $u^* \in U(K)$ .

*Proof:* Consider a control  $u$ . Control  $u'$  such that  $u = u'$  except for  $(i_1, k_1, l_1)$  and

$$\theta_{i_1,j}(k_1, l_1; u') = \begin{cases} 1, & \text{for } j = \underline{a}(k_1, l_1), \\ 0, & \text{for } j \neq \underline{a}(k_1, l_1), \end{cases} \quad (3.7)$$

reduces  $q_k$  by Lemma 2.1, where  $k > i_1 + \underline{a}(k_1, l_1) - 1(i_1 \neq 0) = y(i_1, \underline{a}(k_1, l_1))$ . That is,

$$\begin{cases} q_k(u) \geq q_k(u'), & \text{for } k > y(i_1, \underline{a}(k_1, l_1)), \\ q_k(u) = q_k(u'), & \text{for } k \leq y(i_1, \underline{a}(k_1, l_1)). \end{cases} \quad (3.8)$$

Therefore, if  $y(i_1, \underline{a}(k_1, l_1)) \geq K$ , then  $q_k(u) \geq q_k(u')$  for  $k \geq y(i_1, \underline{a}(k_1, l_1)) \geq K$ , while  $q_k(u) = q_k(u')$  for  $k < K$ .

If  $\theta_{i_1,j_1}(k_1, l_1; u) > 0$  for such  $(i_1, j_1, k_1, l_1)$  that  $y(i_1, j_1) \geq K, y(i_1, \underline{a}(k_1, l_1)) < K \leq y(i_1, \bar{a}(k_1, l_1))$ , then there is a control  $u'$  such that  $u = u'$  except for  $(i_1, k_1, l_1)$ , and

$$\theta_{i_1,j}(k_1, l_1; u') = \begin{cases} = 0, & \text{for } y(i_1, j) \geq K, \\ \geq \theta_{i_1,j}(k_1, l_1; u), & \text{for } y(i_1, j) = K - 1, \\ = \theta_{i_1,j}(k_1, l_1; u), & \text{for } y(i_1, j) < K - 1. \end{cases} \quad (3.9)$$

Using Lemma 2.1,

$$\begin{cases} q_k(u) \geq q_k(u'), & \text{for } k > K - 1, \\ q_k(u) = q_k(u'), & \text{for } k \leq K - 1. \end{cases} \quad (3.10)$$

The above mentioned procedures (3.7)-(3.10) for constructing the controls  $u'$  from the control  $u$ , show that if control  $u \notin U(K)$ , then we can derive such a control  $u'$  that  $q_k(u) \geq q_k(u')$  for some  $k \geq K$ , while  $q_k(u) = q_k(u')$  for all  $k < K$ . In other words, the controls in  $U(K)$  minimize  $q_k$ , where  $k \geq K$ , keeping  $q_k$  ( $k < K$ ) unchanged. Thus, for fixed  $q_1, \dots, q_{K-1}$ , minimizing  $q_k$  ( $k \geq K$ ) means maximizing  $J_c(u) = \sum_{k=1}^{\infty} (c - k)q_k$ , since  $c - k < 0$  for  $k \geq K$ . Therefore,  $u^*$ , the optimal control of (3.4), is in  $U(K)$ . ■

Define  $J(i; u^*) = \sum_{k=i}^{\infty} (c - k)q_k$ . Set  $\xi_{VD}(i) = \sum_{j,k} \xi_D(j) \lambda_V(k) 1(i = \underline{a}(k, l))$ , the probability that  $a_{VD} = i$  under the control minimizing the accepted voice group size.



$J(K; u^*)$  can then be expressed by  $q_1(u^*), \dots, q_{K-1}(u^*)$ , employing Appendix 1.

$$\begin{aligned} J(K; u^*) &= \frac{\sum_{l=K}^{\infty} q_0 \xi_{VD}(l) \{c(K-l-1) - (K(K-1) - l(l+1))\}}{2(\rho_L - 1)} \\ &+ \frac{\sum_{l=K}^{\infty} q_0 \xi_{VD}(l) (K-l-1) F_L''(1)}{2(\rho_L - 1)^2} \\ &+ \frac{\sum_{j=1}^{K-1} \sum_{l=K+1-j}^{\infty} q_j \xi_{VD}(l) \{c(K-l-j) - (K(K-1) - (l+j)(l+j-1))\}}{2(\rho_L - 1)} \\ &+ \frac{\sum_{j=1}^{K-1} \sum_{l=K+1-j}^{\infty} q_j \xi_{VD}(l) (K-l-j) F_L''(1)}{2(\rho_L - 1)^2} \\ &\triangleq \sum_{j=0}^{K-1} g_{K,j} q_j. \end{aligned}$$

Here,  $F_L''(1) = m_{V2} - \lambda_V + \lambda_V N_L(N_L - 1) + 2\lambda_V N_L \lambda_D h + m_{D2} - \lambda_D + \lambda_D(h_2 - h)$ .

$J(K; u^*)$  can be expressed by a linear combination of  $q_1, \dots, q_{K-1}$ . Thus, using the fact that

$$J_c(u^*) = J(K; u^*) + \sum_{k=0}^{K-1} (c-k) q_k,$$

$J_c(u^*)$  can be also expressed by a linear combination of  $q_1, \dots, q_{K-1}$ . Hence,  $J_c(u^*)$  can be expressed by a linear combination of  $q_1, \dots, q_{i-1}$  ( $i \leq K$ ), using (3.6) repeatedly. Denote  $\alpha_{i,k}$  as the weighting coefficient. That is,

$$J_c(u^*) = \sum_{k=0}^{i-1} \alpha_{i,k} q_k.$$

Set  $\alpha_{i-1} = \alpha_{i,i-1}$  to simplify notation. Then,

$$J_c(u^*) = \alpha_{i-1} q_{i-1} + \sum_{k=0}^{i-2} \alpha_{i,k} q_k. \quad (3.11)$$

Roughly speaking, for fixed  $q_k$  ( $k = 0, \dots, i-2$ ), if  $\alpha_{i-1} < 0$ , then minimizing  $q_{i-1}$  is optimal; if  $\alpha_{i-1} \geq 0$ , then maximizing  $q_0, \dots, q_{i-1}$  is optimal. Formal results are stated in the following theorem.

$\{\alpha_i\}$  can be provided as

$$\alpha_i = c - i + (1 - \lambda)^{-1} \{\alpha_{i+1} \eta_{i,2} + \alpha_{i+2} \eta_{i,3} + \dots + \alpha_{K-1} \eta_{i,K-i}\} + g_{K,i}, \quad (K \geq i \geq 1), \quad (3.12)$$

$$\alpha_0 = c + (1 - \lambda)^{-1} \{\alpha_1 \eta_{0,1} + \dots + \alpha_{K-1} \eta_{0,K-1}\} + g_{K,0}. \quad (3.13)$$

Then we define the randomized stationary control  $v_n(\tau)$  for each  $n$ , which at each decision epoch employs either  $u_n$  or  $u_{n+1}$ , with probability  $\tau$  and  $1 - \tau$ . Here,  $u_n$  is defined by Eq. (3.5).

**Theorem 1.** There exist an integer  $n$  and  $\tau \in [0, 1]$  such that  $v_n(\tau)$  is the optimal control of (3.4).

*Proof:* There exists an integer  $K$  such that  $K-1 < c \leq K$ . Then, employing Lemma 3,  $u^* \in U(K)$ . Consider a control  $u \in U(K)$ . As in the proof of Lemma 3, we can construct control  $u' \in U(K-1) \subset U(K)$  for the control  $u$  such that  $q_{K-1}(u) \geq q_{K-1}(u')$  and  $q_k(u) = q_k(u')$  ( $k < K-1$ ). Then, if  $\alpha_{K-1} < 0$ ,  $J_c(u) \leq J_c(u')$ . That is,  $u^* \in U(K-1)$ .

If  $\alpha_{K-2} < 0$ , then  $u^* \in U(K-2)$ , if  $\alpha_{K-3} < 0$ , then  $u^* \in U(K-3)$ , and so on.

Assume that a control  $u \in U(i)$ . Consider a control  $u' \in U(i)$  which is equal to  $u$  except for  $(i_1, k_1, l_1)$  and

$$\theta_{i,j}(k_1, l_1; u') = \begin{cases} 1, & \text{for } j = j_1, \\ 0, & \text{for } j \neq j_1. \end{cases}$$

Here,  $y(i_1, j_1) = i-1$ ,  $y(i_1, \underline{a}(k_1, l_1)) < i \leq y(i_1, \bar{a}(k_1, l_1))$ . Since  $u \in U(i)$  and  $y(i_1, j_1) = i-1$ ,  $\theta_{i,j}(k_1, l_1; u) = 0$  for  $j > j_1$ . Employing Lemma 2.2,  $q_k(u) \leq q_k(u')$  for all  $k$ . It is shown in Appendix 2 that  $\alpha_{i,k} > 0$  ( $k = 0, 1, 2, \dots, i-2$ ), if  $\alpha_{i-1} \geq 0$ . Therefore,  $u^* = u_i = v_{i-1}(0) = v_i(1)$ , if  $\alpha_{K-1} < 0, \alpha_{K-2} < 0, \dots, \alpha_i < 0$  and  $\alpha_{i-1} > 0$ . If



$\alpha_{K-1} < 0, \alpha_{K-2} < 0, \dots, \alpha_i < 0, \alpha_{i-1} = 0$ , then  $u^* = v_{i-1}(\tau)$ , for all  $\tau$ . If  $\alpha_{K-1} < 0, \alpha_{K-2} < 0, \dots, \alpha_0 < 0$ , then  $u^* \in U(0)$ . That is,  $u^* = u_0 = v_0(1)$ . This concludes the proof. ■

**Lemma 4.** The constraint in (3.2),  $\tilde{C}(u = v_i(\tau))$ , is a continuous and decreasing function of  $\tau$  over the interval  $[0, 1]$  for each  $i$ .

*Proof:* Under the control  $v_i(\tau)$ , it is possible to obtain

$$G(z) \triangleq \sum_{i=0}^{\infty} q_i z^i = \frac{q_0(F_L(z) - zF_0(z)) + \sum_{j=1}^{i-1} q_j z^j (F_L(z) - F_j(z))}{F_L(z) - z}. \quad (3.14)$$

Here,

$$\begin{aligned} F_j(z) &= \sum_{k,l,m} \theta_{j,m}(k, l; u) \lambda_V(k) \xi_D(l) z^m \\ &= \sum_{k,l,m} \lambda_V(k) \xi_D(l) z^m \\ &\quad \{1(y(j, \underline{a}(k, l)) > i, m = \underline{a}(k, l)) \\ &\quad + 1(y(j, \bar{a}(k, l)) \leq i, m = \bar{a}(k, l)) \\ &\quad + 1(y(j, \underline{a}(k, l)) \leq i < y(j, \bar{a}(k, l))) \\ &\quad (\tau 1(i = y(j, m)) + (1 - \tau) 1(i + 1 = y(j, m)))\}. \end{aligned} \quad (3.15)$$

Therefore,

$$\begin{aligned} G'(1) &= \frac{\{q_0(F_L''(1) - 2F_0'(1) - F_0''(1)) + 2 \sum_{j=1}^{i-1} q_j j(\rho_L - F_j'(1)) + \sum_{j=1}^{i-1} q_j (F_L''(1) - F_j''(1))\}}{2(\rho_L - 1)} \\ &\quad - \frac{\{q_0(\rho_L - 1 - F_0'(1)) + \sum_{j=1}^{i-1} q_j (\rho_L - F_j'(1))\} F_L''(1)}{2(\rho_L - 1)^2}, \end{aligned} \quad (3.16)$$

and

$$G(1) = \frac{\{q_0(\rho_L - 1 - F_0'(1)) + \sum_{j=1}^{i-1} q_j (\rho_L - F_j'(1))\}}{\rho_L - 1}. \quad (3.17)$$

Consequently,  $\tilde{C}(u = v_i(\tau)) = G'(1) - (1 + Q)(G(1) - 1)$  is a continuous and decreasing function of  $\tau$  for each  $i$ . ■

**Lemma 5.** For any  $i \geq 0$  and  $\tau \in [0, 1]$ , there exists a  $c$  such that  $v_i(\tau)$  is an optimal control for the criterion  $J_c$ .

*Proof:* We show by mathematical induction that for any  $i \geq 0$  and  $\tau \in [0, 1]$ ,  $v_i(\tau)$  can be an optimal control, when  $c = 0 \rightarrow \infty$ . Roughly speaking,  $u_i$  is an optimal control and the optimal  $i$  increases by units of one, as  $c$  increases continuously. In other words, if  $u^* = u_i$  at  $c = c_0$ , then  $u^* = u_i$  or  $u_{i+1}$  at  $c = c_0 + 0$ .

(Step 0)

When  $Q \rightarrow \infty$ , then  $c \rightarrow \infty$  and the original problem (2.1)-(2.2) is reduced to an optimization problem without a constraint. Therefore, the optimal control is  $u_\infty$ .

(Step 1)

When  $c = 0$ ,  $u_0 = v_0(1)$  is optimal, by Lemma 3.

(Step 2)

Assume that  $K - 1 < c \leq K$  and that  $\alpha_{K-1} < 0, \dots, \alpha_i < 0, \alpha_{i-1} \geq 0$ . Then  $u_i = v_i(1) = v_{i-1}(0)$  is optimal.  $\alpha_{K-1} = c - K + g_{K,K-1}$  is a continuous and increasing function of  $c$ . Therefore,  $\alpha_{K-2}$  is also a continuous and increasing function of  $c$  by (3.12).  $\alpha_{K-3}$  is a continuous and increasing function of  $c$ , because  $\alpha_{K-1}$  and  $\alpha_{K-2}$  are continuous and increasing functions. Consequently,  $\alpha_{K-1} < 0, \dots, \alpha_i < 0$  is a continuous and increasing function of  $c$ .

Increase  $c$ , while  $c \leq K$ . Then,  $\alpha_{K-1} < 0, \dots, \alpha_i < 0$  also increase.

(Step 2a)



If  $\alpha_j (i < j \leq K-1)$  approaches 0, it is shown using Appendix 3 that  $\alpha_{j-1} > \alpha_j + 1$ . Therefore,  $\alpha_{j-1} > 0$ . Consequently, the first among  $\alpha_{K-1}, \dots, \alpha_i$  which becomes positive when  $c$  increases is  $\alpha_i$ . That is,  $\alpha_{K-1} < 0, \dots, \alpha_{i+1} < 0$  and  $\alpha_i > 0$  for some value of  $c$ . This means  $u_{i+1} = v_{i+1}(1) = v_i(0)$  is optimal. In addition, there exists a  $c$  such that  $\alpha_i = 0$ , because of the continuity of  $\alpha_i$ . Then,  $u^* = v_i(r)$  for all  $r \in [0, 1]$ .

(Step 2b)

Consider the case that  $\alpha_{K-1} < 0, \dots, \alpha_i < 0, \alpha_{i-1} \geq 0$  at  $c = K$ . At  $c = K + 0$ ,

$$\alpha_i(K+0) = c - i + (1 - \lambda)^{-1} \{ \alpha_{i+1}(K+0)\eta_{i,2} + \dots + \alpha_K(K+0)\eta_{i,K-i+1} \} + g_{K+1,i}. \quad (3.18)$$

Here,  $\alpha_i(K+0)$  shows the dependence of  $\alpha_i$  on  $c (= K+0)$ . At  $c = K$ ,

$$J(K; u) = \sum_{k=K}^{\infty} (c - k)q_k = \sum_{k=K+1}^{\infty} (c - k)q_k = J(K+1; u)$$

for all  $u$ . Since  $J(K; u) = \sum_{j=0}^{K-1} g_{K,j}q_j$  and  $J(K+1; u) = \sum_{j=0}^K g_{K+1,j}q_j$ , employing (3.12),

$$g_{K,i} = g_{K+1,i} + g_{K+1,K}(1 - \lambda)^{-1}\eta_{i,K-i+1}$$

$$= g_{K+1,i} + \alpha_K(K+0)(1 - \lambda)^{-1}\eta_{i,K-i+1},$$

for  $0 \leq i \leq K-1$ . Using the equation above, (3.12) and (3.18), we obtain

$$\alpha_i(K+0) = \alpha_i(K).$$

Therefore,  $\alpha_{K-1} < 0, \dots, \alpha_i < 0, \alpha_{i-1} \geq 0$  at  $c = K+0$ , when  $\alpha_{K-1} < 0, \dots, \alpha_i < 0, \alpha_{i-1} \geq 0$  at  $c = K$ .

Using (Step 2) for  $K < c \leq K+1, K+1 < c \leq K+2$  and so on,  $u_{i+1} = v_{i+1}(1) = v_i(0)$  and  $v_i(r)$  can be an optimal control.

This concludes the proof. ■

We are now in a position to state the main theorem.

**Theorem 2.** There exists an optimal control  $u^*$  of (2.1) with a constraint (2.2) in  $V_r \triangleq \bigcup_i V_r(i)$  if  $\bar{Q}_0 \leq Q \leq \bar{Q}_\infty$ . Here,  $V_r(i)$  is the set of controls  $\{v_i(r) \mid r \in [0, 1]\}$ . There is no optimal control that meets the constraint if  $\bar{Q}_0 > Q$ . The deterministic control  $u_\infty$  is optimal if  $\bar{Q}_\infty \leq Q$ . Here,  $\bar{Q}_0$  and  $\bar{Q}_\infty$  are the average queue lengths under the controls  $u_0$  and  $u_\infty$ ;

$$\bar{Q}_0 = \frac{((m_{V2} - \lambda_V)N_L^2 + \lambda_V N_L(N_L - 1) + 2\lambda_V N_L \lambda_D h + (m_{D2} - \lambda_D)h^2 + \lambda_D(h_2 - h))}{2(1 - \rho_L)},$$

$$\bar{Q}_\infty = \frac{((m_{V2} - \lambda_V)N_V^2 + \lambda_V N_V(N_V - 1) + 2\lambda_V N_V \lambda_D h + (m_{D2} - \lambda_D)h^2 + \lambda_D(h_2 - h))}{2(1 - \rho_H)},$$

where  $\rho_H = \lambda_D h + \lambda_V N_V$ .

*Proof:* First assume  $\bar{Q}_0 \leq Q \leq \bar{Q}_\infty$ . Then, there exists an integer  $n$  such that  $C(u_n) \leq Q < C(u_{n+1})$ . With Lemma 4, it is shown that there exists an  $r \in [0, 1]$  such that  $\tilde{C}(v_n(r)) = Q$ . Employing Lemma 5, there exists a  $c$  such that the optimal control of (3.4) is  $v_n(r)$ . Thus, for any control  $u$ ,

$$\tilde{J}(v_n(r)) - \gamma \tilde{C}(v_n(r))$$

$$= J_\gamma(v_n(r))$$

$$\geq J_\gamma(u)$$

$$= \tilde{J}(u) - \gamma \tilde{C}(u).$$

Therefore,

$$\tilde{J}(v_n(r)) - \tilde{J}(u)$$

$$\geq \gamma(\tilde{C}(v_n(r)) - \tilde{C}(u))$$

$$= \gamma(Q - \tilde{C}(u)) > 0.$$



Consequently, it is shown with Lemma 1 that  $v_n(r)$  is an optimal control of (2.1) which meets the constraint (2.2). For other cases, it can be observed that  $J(v_n(r))$  is an increasing function of  $n$  and a decreasing function of  $r$ , and that  $C(v_n(r))$  is a decreasing function of  $n$  and an increasing function of  $r$ . In addition,  $\bar{Q}_0$  and  $\bar{Q}_\infty$  are the average queue lengths of cells using controls  $u_0$  and  $u_\infty$ . Now if  $\bar{Q}_0 > Q$ , the average queue length using any control is greater than  $Q$ , and if  $\bar{Q}_\infty \leq Q$ , control  $u_\infty$  maximizes  $E[N]$ . ■

#### 4. The optimal parameters

Results in the preceding section provides the following procedure to obtain the optimal parameters  $n$  and  $r$  for the constrained optimal control  $v_n(r)$ .

(Step 1)

Find an integer  $n^*$  such that  $C(u_{n^*}) \leq Q < C(u_{n^*+1})$ , if  $\bar{Q}_0 \leq Q \leq \bar{Q}_\infty$ . Here,  $C(u_n) = C(v_{n-1}(0))$  can be obtained from the following equations with  $r = 0$ .

$$C(v_{n-1}(r)) = G'(1) * p_0 - 1 + p_0. \quad (4.1)$$

$$p_0 = (G(1))^{-1}. \quad (4.2)$$

$G(1)$  and  $G'(1)$  are provided in (3.16) and (3.17). Using the fact that

$$C(u_0) < C(u_1) < \dots$$

makes it easy to find  $n^*$ .

If  $\bar{Q}_0 > Q$ , there is no optimal control which meets the constraint.

If  $\bar{Q}_\infty \leq Q$ ,  $u_\infty$  is optimal.

(Step 2)

Find an  $r^* \in [0, 1]$  such that  $C(v_{n^*}(r^*)) = Q$  using (4.1) and (4.2). Then,  $v_{n^*}(r^*)$  is an optimal control of (2.1) which meets the constraint (2.2). The property shown in Lemma 4 that the constraint  $C$  is a decreasing and continuous function of  $r$ , is useful for obtaining  $r^*$ .

#### 5. An average waiting time constraint

An average waiting time constraint used instead of an average queue length constraint is considered.

$$C_w(u) \triangleq E[\text{cell waiting time}] \leq W \quad (5.1)$$

By Little's formula,

$$C_w(u) = (\lambda_V E[N] + \lambda_D h)^{-1} \sum_{k=1}^{\infty} (k-1)p_k \leq W.$$

Thus, using  $\lambda_V E[N] + \lambda_D h = 1 - p_0$ ,

$$\sum_{k=1}^{\infty} (k-1)p_k \leq W(1 - p_0).$$

Introducing  $\{q_k\}$  and using  $p_0 = (1 + \sum_{k=1}^{\infty} q_k)^{-1}$ ,

$$\sum_{k=1}^{\infty} (k-1-W)q_k \leq 0.$$

Therefore, our problem is equivalent to

$$\begin{aligned} & \underset{u}{\text{maximize}} && \sum_{k=1}^{\infty} q_k, \\ & \text{subj. to} && \sum_{k=1}^{\infty} (k-1-W)q_k \leq 0. \end{aligned} \quad (5.2)$$

$$\quad (5.3)$$



Equations (5.2) and (5.3) have the same structure as (3.1) and (3.2). The results in Section 3 were derived without using the fact that the right hand side of (3.2) =  $Q$ . Therefore, the argument in the previous section is also valid for the problem stated in (5.2)-(5.3). In other words, the optimal control with an average waiting time constraint is also  $v_i(r)$ , which is a randomized modification of two deterministic controls with saturation  $u_i, u_{i+1}$ .

## 6. Conclusions

In this chapter, it was proved that the optimal control for maximizing an average voice coding rate under an average queue length or an average waiting time constraint is randomization of two deterministic controls which have simple structure (feedback with input saturation) and that the control can be specified by two parameters. The intuitive control scheme for reducing bit rates which has been analyzed and evaluated in the literature and implemented in experimental systems, is the bang-bang control type. Therefore, our results require a slight modification for that scheme.

An infinite buffer system was considered. However, this analysis can also be extended to finite buffer systems, for which it can be shown that the optimal control has the same structure. For finite buffer systems, it is known that attention can be restricted to the class of randomized stationary controls even if we have interest for the class of controls which depend on the past history of the system [Derman 70]. Therefore, this optimal control for finite buffer systems is also optimal in the class of controls that may depend on the past history of the system.

Arrival processes were assumed to be independent among slots. The same structure, randomization of two feedback controls with saturation, is shown to be optimal also for

highly correlated arrival processes [Saito 90b], using the results and methodology presented here. It is shown that the optimality is also valid for the constraint on the quantile of voice cell waiting time.

Additionally, our results are valid for other queueing disciplines as long as they are work conserving.



## Appendix 1

Define  $G_K(z) = \sum_{k=K}^{\infty} q_k z^k$ . Under the control  $u \in U(K)$ ,  $\eta_{j,k} = \sum_{l=k}^{\infty} \xi_{VD}(l)$  for  $g_{j,k} \geq K$ . Employing (3.6), we obtain

$$G_K(z) = \frac{\sum_{l=K}^{\infty} \xi_{VD}(l) q_0 (z^K - z^{l+1}) + \sum_{j=1}^{K-1} \sum_{l=K+1-j}^{\infty} \xi_{VD}(l) q_j (z^K - z^{l+j})}{F_L(z) - z}$$

Here,  $F_L(z)$  denotes the generating function of the number of accepted cells provided that control is used to reduce the group size of the arriving voice cells to  $N_L$ . Since, the generating functions of the number of groups of arriving voice cells and data cells are  $M_V(z)$  and  $M_D(z)$ , and the generating function of the number of cells in a data group is  $H(z)$ ,

$$F_L(z) = \sum_{l=0}^{\infty} \xi_{VD}(l) z^l = M_V(z^{N_L}) M_D(H(z)).$$

Thus, we obtain

$$\begin{aligned} G_K(1) &= \frac{\sum_{l=K}^{\infty} q_0 \xi_{VD}(l) (K-l-1) + \sum_{j=1}^{K-1} \sum_{l=K+1-j}^{\infty} q_j \xi_{VD}(l) (K-l-j)}{2(\rho_L - 1)}, \\ G'_K(1) &= \frac{\sum_{l=K}^{\infty} q_0 \xi_{VD}(l) \{ (K(K-1) - l(l+1)) \}}{2(\rho_L - 1)} \\ &\quad - \frac{\sum_{l=K}^{\infty} q_0 \xi_{VD}(l) (K-l-1) F'_L(1)}{2(\rho_L - 1)^2} \\ &\quad + \frac{\sum_{j=1}^{K-1} \sum_{l=K+1-j}^{\infty} q_j \xi_{VD}(l) \{ K(K-1) - (l+j)(l+j-1) \}}{2(\rho_L - 1)} \\ &\quad - \frac{\sum_{j=1}^{K-1} \sum_{l=K+1-j}^{\infty} q_j \xi_{VD}(l) (K-l-j) F'_L(1)}{2(\rho_L - 1)^2}. \end{aligned}$$

## Appendix 2

Assume that  $\alpha_{K-1}, \dots, \alpha_i < 0$  and  $\alpha_{i-1} \geq 0$ . Then, it is shown in the following that  $\alpha_{i,k} > 0$  for  $0 \leq k < i-1$ . Notice that  $\sum_{m=l-k}^{\infty} \xi_{VD}(m) = \eta_{k,l-k} > \eta_{j,l-j} = \sum_{m=l-j}^{\infty} \xi_{VD}(m)$ ,  $\eta_{k,l-k} > \eta_{0,l-1} = \sum_{m=l-1}^{\infty} \xi_{VD}(m)$  for  $0 \leq j < k \leq i-1$  and  $i+1 \leq l \leq K$ , since  $u^* \in U(i)$ . Additionally, it can be shown that  $g_{K,K-1} < \dots < g_{K,2} < g_{K,1} = g_{K,0} < 0$ .

$$\alpha_{i,k} - \alpha_{i,j} = j - k + g_{K,k} - g_{K,j}$$

$$+ (1-\lambda)^{-1} \{ \alpha_i (\eta_{k,i+1-k} - \eta_{j,i+1-j}) + \dots + \alpha_{K-1} (\eta_{k,K-k} - \eta_{j,K-j}) \} < 0$$

$$\alpha_{i,k} - \alpha_{i,0} = -k + g_{K,k} - g_{K,0}$$

$$+ (1-\lambda)^{-1} \{ \alpha_i (\eta_{k,i+1-k} - \eta_{0,i}) + \dots + \alpha_{K-1} (\eta_{k,K-k} - \eta_{0,K-1}) \} < 0$$

Therefore,  $\alpha_{i,0} > \alpha_{i,1} > \dots > \alpha_{i,i-1} = \alpha_{i-1} \geq 0$ .

## Appendix 3

Assume that  $\alpha_{K-1} < 0, \dots, \alpha_i < 0, \alpha_{i-1} \geq 0$  and  $K-1 < c \leq K$ . Notice that  $\eta_{j,l} > \eta_{j-1,l+1}$  for  $i < j \leq K-1$  and  $l > 0$ , since  $u^* \in U(i)$ . It can be shown that  $g_{K,K-1} < \dots < g_{K,2} < g_{K,1} = g_{K,0} < 0$ . Therefore, when  $\alpha_j \uparrow 0$ ,

$$\alpha_j - \alpha_{j-1}$$

$$= -1 + g_{K,j} - g_{K,j-1} - (1-\lambda)^{-1} \alpha_j \eta_{j-1,2}$$

$$+ (1-\lambda)^{-1} \{ \alpha_{j+1} (\eta_{j,2} - \eta_{j-1,3}) + \dots + \alpha_{K-1} (\eta_{j,K-j} - \eta_{j-1,K-j+1}) \} < -1.$$



#### Appendix 4

$N_V$  = # of cells in a group of voice cells (constant).

$m_V$  = # of groups of voice cells arriving in a slot.

$$\lambda_V(k) = \Pr(m_V = k).$$

$$\lambda_V = E[m_V].$$

$$m_{V2} = E[m_V^2].$$

$$M_V(z) = \sum_{k=0}^{\infty} \lambda_V(k) z^k.$$

$N_D$  = # of cells in a group of data cells.

$m_D$  = # of data cell groups arriving in a slot.

$$\lambda_D(k) = \Pr(m_D = k).$$

$$\lambda_D = E[m_D].$$

$$m_{D2} = E[m_D^2].$$

$$M_D(z) = \sum_{k=0}^{\infty} \lambda_D(k) z^k.$$

$$h(i) = \Pr(N_D = i).$$

$$h = E[N_D].$$

$$H(z) = \sum_{k=0}^{\infty} h(k) z^k.$$

$a_D$  = Total number of data cells arriving in a slot.

$$\xi_D(k) = \Pr(a_D = k).$$

$N$  = # of cells accepted in a group of voice cell ( $N_L \leq N \leq N_V$ ).

$$\rho_L = \lambda_D h + \lambda_V N_L.$$

$$\rho_H = \lambda_D h + \lambda_V N_V.$$

$z$  = Total number of cells in the system.

$$p_k = \Pr(z = k).$$

$$q_k = p_k / p_0.$$

$Q$  = Queue length constraint (See (2.2)).

$a_{VD}$  = Total number of accepted voice and data cells.

$$\theta_{i,j}(k, l; u) = \Pr(a_{VD} = j \mid z = i, m_V = k, a_D = l, \text{control} = u).$$

$\underline{a}(i, j)$  = Minimum number of accepted cells =  $iN_L + j$ .

$\bar{a}(i, j)$  = Maximum number of accepted cells =  $iN_V + j$ .

$$\eta_{i,j} = \Pr(a_{VD} \geq j \mid z = i) = \sum_{m=j}^{\infty} \sum_{k=0}^{\infty} \sum_{l=0}^{\infty} \theta_{i,m}(k, l; u) \lambda_V(k) \xi_D(l).$$

$y(i, j)$  = Total number of cells in the system after accepting  $j$  data and voice cells,

when  $z = i$ .

$$= i + j - 1 (i \neq 0).$$

$$\xi_{VD}(i) = \Pr(a_{VD} = i \mid N = N_L).$$



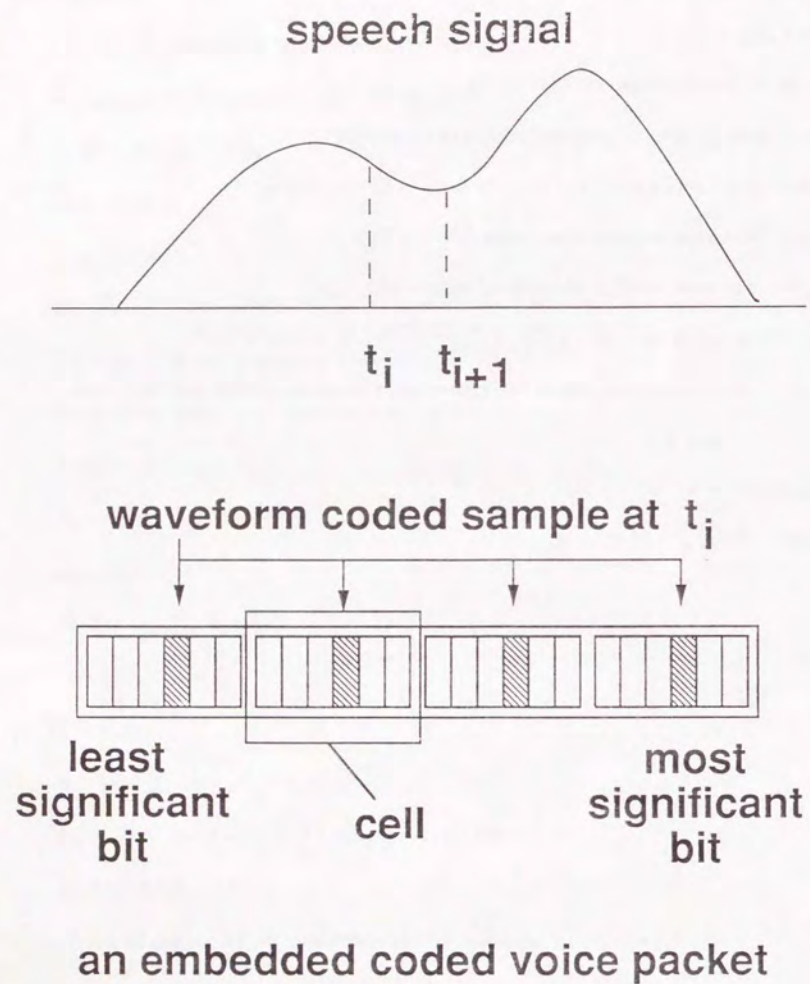


Figure 1. Embedded coding

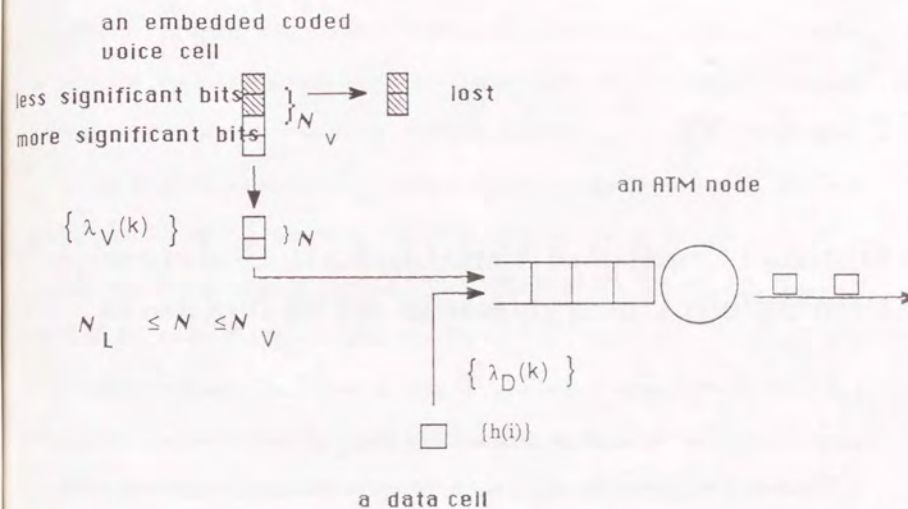


Figure 2. Model