

Quality of Service over ATM Networks

by

Steven B. Winstanley

SUBMITTED FOR THE DEGREE OF DOCTOR OF PHILOSOPHY

Department of Electronic Engineering

Queen Mary and Westfield College

University of London

1998

Abstract

The thesis examines Quality of Service (QoS) over Asynchronous Transfer Mode (ATM) networks. QoS over traditional communication networks has generally been taken for granted, but many newer networks, such as the Frame Relay and Fast Ethernet, fail to supply suitable stringent guarantees. The definitions of network performance and QoS are identified; the protocol stack and the control function in terms of network performance and QoS are dissected. This enables mapping strategies, points of demarcation, and measurement methods to be assessed and combinations of these are manipulated to provide possible solutions for optimal ATM network operation.

Quality of service has been defined as a property measured between the physical and the application layer on the International Standards Organisation (ISO) seven layer model. A user requires application-layer to application-layer QoS, but, the network operator can only guarantee the ATM **resources** required for QoS, and cannot provide true end-to-end QoS guarantees. In this thesis, a flexible QoS framework has been proposed that is simple to use, beneficial to customers and operators, and flexible enough to facilitate both circuit switched and IP services.

The research in this thesis covers the real-time performance of ATM networks, the feasibility of high quality CBR services, QoS to network performance mapping, the user benefit of shaping while maintaining QoS, and the increased network utilisation by incorporating end-system shaping. The new approach abolishes the rt-VBR category and in so doing, removes the buffering requirements for real-time circuits from the network to the end-systems.

The research on CBR ATCs enables bounded requests for network resources to be made to the ATM layer in order to provide sufficient QoS to CBR services. When Variable Bit Rate (VBR) traffic is introduced to ATM networks new problems arise. VBR traffic can waste expensive network resources and customers have difficulty specifying complex traffic contracts. VBR traffic can be moulded into CBR traffic descriptors that the customer specifies by using shaping devices. Research here has been on reducing the reserved network resources while

maintaining a high QoS. The author's research has shown that the customer and network operator make benefits by reducing the application's reserved capacity and obtaining greater multiplexing densities in switch queues.

Using the CBR and VBR results from the research, this thesis provides new recommendations that can be made to users and operators for a new QoS framework.

Table of Contents

1. INTRODUCTION	1
1.1 OVERVIEW	1
1.2 THE ATM TESTBED.....	5
1.3 SUMMARY OF CONTRIBUTION	6
1.4 THESIS LAYOUT	7
2. THE ATM PROTOCOL STACK	8
2.1 THE APPLICATION LAYER	8
2.2 THE AAL LAYER.....	9
2.2.1 <i>ATM Adaptation Layer 1</i>	11
2.2.2 <i>ATM Adaptation Layer 2</i>	12
2.2.3 <i>ATM Adaptation Layer 3/4</i>	12
2.2.4 <i>ATM Adaptation Layer 5</i>	13
2.3 THE ATM LAYER	13
2.3.1 <i>Constant Bit Rate (CBR)</i>	14
2.3.2 <i>Real-time Variable Bit Rate (rt-VBR)</i>	14
2.3.3 <i>Non-real-time Variable Bit Rate (nrt-VBR)</i>	15
2.3.4 <i>Available Bit Rate (ABR)</i>	15
2.3.5 <i>Unspecified Bit Rate (UBR)</i>	16
2.3.6 <i>ATM Block Transfer (ABT)</i>	17
2.4 TRAFFIC CONTROL FUNCTIONS	18
2.4.1 <i>Connection Admission Control</i>	18
2.4.2 <i>Feedback</i>	18
2.4.3 <i>Policing</i>	18
2.4.4 <i>Priority</i>	19
2.4.5 <i>Network Resource Management</i>	19
2.4.6 <i>Frame or Packet Discard</i>	19
2.4.7 <i>ABR Flow Control</i>	20
2.5 ITU-T AND ATMF QoS MAPPING.	20
2.6 SERVICE CATEGORIES AND QUALITY OF SERVICE CLASSES	21

2.6.1	<i>Service Categories</i>	23
2.6.2	<i>Definition of Service Categories with Network Performance Parameters</i>	23
2.7	APPLICATION CATEGORIES, AALS AND ATM SERVICE CATEGORIES.....	25
2.7.1	<i>Overview</i>	25
2.7.2	<i>Functionality Mapping Requirements</i>	26
2.7.3	<i>Quality of Service Mapping Requirements</i>	28
2.7.4	<i>Mapping Pricing and Traffic Control</i>	32
2.8	SUMMARY.....	33
3.	ATM PERFORMANCE PARAMETERS	35
3.1	NETWORK PERFORMANCE GUARANTEES.....	35
3.2	CELL EVENTS AND MEASUREMENT POINTS.....	36
3.3	CELL TRANSFER OUTCOME.....	38
3.4	NEGOTIATED NETWORK PERFORMANCE PARAMETERS.....	38
3.4.1	<i>Information Transfer Speed</i>	39
3.4.2	<i>Information Transfer Dependability</i>	45
3.5	NON-NEGOTIATED NETWORK PERFORMANCE PARAMETERS.....	47
3.5.1	<i>Information Transfer Accuracy</i>	47
3.6	SOURCES OF PERFORMANCE DEGRADATION.....	49
3.7	SUMMARY.....	50
4.	PROPOSED QOS FRAMEWORK	51
4.1	ATM NETWORK PERFORMANCE PARAMETERS.....	51
4.2	QUALITY OF SERVICE.....	51
4.3	CUSTOMER, OPERATOR AND SERVICE PROVIDER.....	53
4.4	THE NEW APPROACH.....	54
4.4.1	<i>The Problem</i>	54
4.4.2	<i>The Solution</i>	58
4.4.3	<i>The Research Methodology</i>	61
5.	AN EXPERIMENTAL APPROACH TO ATM QOS	64
5.1	EXPERIMENT 1, CELL TRANSFER DELAY AND CELL DELAY VARIATION MEASUREMENTS.....	64
5.1.1	<i>Introduction</i>	64

5.1.2	<i>Experiment Set-up</i>	65
5.1.3	<i>Results</i>	67
5.1.4	<i>Summary</i>	73
5.2	EXPERIMENT 2, INVESTIGATION OF END-SYSTEM SPEED DEGRADATION BY CBR PLAY-OUT QUEUEING.....	74
5.2.1	<i>Introduction</i>	74
5.2.2	<i>Play-Out Buffering</i>	75
5.2.3	<i>Experimental Set-up</i>	77
5.2.4	<i>Results</i>	78
5.2.5	<i>Summary</i>	81
5.3	EXPERIMENT 3, MAPPING SUBJECTIVE QoS INTO NETWORK PERFORMANCE PARAMETERS.....	82
5.3.1	<i>Introduction</i>	82
5.3.2	<i>Application Description</i>	83
5.3.3	<i>Results</i>	87
5.3.4	<i>Summary</i>	91
5.4	EXPERIMENT 4, PCR REDUCTION BY TRAFFIC SHAPING WHILE MAINTAINING QoS.....	92
5.4.1	<i>Introduction</i>	92
5.4.2	<i>Experiment Set-up</i>	95
5.4.3	<i>Results</i>	97
5.4.4	<i>Summary</i>	103
5.5	EXPERIMENT 5, INCREASING MULTIPLEXING DENSITY BY TRAFFIC SHAPING.....	104
5.5.1	<i>Introduction</i>	104
5.5.2	<i>Experiment Set-up</i>	105
5.5.3	<i>Results</i>	107
5.5.4	<i>Summary</i>	110
6.	DISCUSSION	111
6.1	CELL DELAY.....	112
6.2	EFFECTS OF CDV ON CBR TRAFFIC.....	114
6.3	APPLICATION SUBJECTIVE QoS.....	115
6.4	USER BENEFITS OF SHAPING RT-VBR.....	117

6.5 OPERATOR BENEFITS OF SHAPING RT-VBR	118
6.6 FURTHER WORK	119
7. CONCLUSION.....	121
8. AUTHOR'S PUBLICATIONS	I
9. REFERENCES.....	II

List of Tables

TABLE 2.1: APPLICATION WITH QoS CLASSES	8
TABLE 2.2: QoS CLASSES MAPPED TO APPLICATIONS, AS RECOMMENDED BY THE ATMF	21
TABLE 2.3: QoS CLASSES MAPPED TO ATCs AS RECOMMENDED BY THE ITU-T	22
TABLE 2.4: MAPPING ATCs BETWEEN THE ITU-T AND ATM-FORUM	23
TABLE 2.5: ATM ATTRIBUTES AND SERVICE CATEGORIES SPECIFICATIONS	25
TABLE 2.6: GENERIC PERFORMANCE PARAMETERS FOR QoS ASSESSMENT	28
TABLE 2.7: ATM LAYER QoS PARAMETERS AND GENERIC ASSESSMENT CRITERIA.	29
TABLE 2.8: ATM LAYER QoS CLASSES PROPOSED BY [I.365]	31
TABLE 3.1: NUMBER OF SWITCHES ASSUMED IN EACH PORTION	42
TABLE 3.2: CDV TOLERANCE FOR DIFFERENT PORTIONS	44
TABLE 3.3: DEGRADATION OF NETWORK PERFORMANCE.	49
TABLE 4.1: SIX SIMULATED WORST-CASE TRAFFIC TYPES.....	57
TABLE 4.2: TRAFFIC CHARACTERISTICS AND NETWORK PERFORMANCE COUNTERS.	58
TABLE 5.1: THE ATCs WITH THEIR RESPECTIVE PERFORMANCE PARAMETERS	64
TABLE 5.2: SWITCH DELAY MEASUREMENTS	68
TABLE 5.3: ESTIMATION OF DELAY IN AN EXAMPLE ATM NETWORK.	69
TABLE 5.4: EXPECTED DELAY AND CDV MEASUREMENTS	70
TABLE 5.5: EXPECTED DELAY AND CDV MEASUREMENTS	72
TABLE 5.6: MINIMUM QoS TO CLR MAPPING.....	88
TABLE 5.7: MINIMUM QoS TO CDV MAPPING	90
TABLE 5.8: NETWORK UTILISATION USING ISABEL MULTIMEDIA TERMINALS	109

List of Figures

FIGURE 2.1: THE ATM PROTOCOL STACK AND QoS CLASSES	10
FIGURE 3.1: ATM NETWORK PERFORMANCE REFERENCE POINTS.	37
FIGURE 3.2: ILLUSTRATIVE CELL TRANSFER DELAY PROBABILITY	40
FIGURE 4.1: THE DUAL BUCKET POLICING PARAMETERS	56
FIGURE 4.2: WORST CASE TRAFFIC FROM A DUAL LEAKY BUCKET MECHANISM ...	57
FIGURE 4.3: A SWITCH MULTIPLEXING RT-VBR SOURCES.....	58
FIGURE 5.1: HARDWARE CONFIGURATION TO TEST A SINGLE NETWORK ELEMENT...65	
FIGURE 5.2: BASEL-LEIDSCHENDAM-BASEL PORTIONED	67
FIGURE 5.3: BASEL-OTTAWA-BASEL PORTIONED	67
FIGURE 5.4: EXAMPLE NETWORK.	69
FIGURE 5.5: CELL TRANSFER DELAY BASEL-LEIDSCHENDAM-BASEL.....	71
FIGURE 5.6: CELL TRANSFER DELAY BASEL-OTTAWA-BASEL.....	73
FIGURE 5.7: RECONSTRUCTING THE TRAFFIC PROFILE WITH A PLAY-OUT BUFFER	76
FIGURE 5.8: PLAY-OUT QUEUE EXPERIMENT SET-UP.....	78
FIGURE 5.9: QUEUE OCCUPANCY FOR 0.210 & 38.88 MBIT/S CBR SOURCES.....	79
FIGURE 5.10: MAPPING σ_p OF THE ARRIVAL PROCESS TO σ_N OF THE QUEUE P.D.F...80	
FIGURE 5.11: DELAY WITHIN A CBR END-TERMINAL.....	81
FIGURE 5.12: EXPERIMENTAL SET-UP	83
FIGURE 5.13: 64 KBIT/S N-ISDN TELEPHONY.....	84
FIGURE 5.14: 384 KBIT/S HQ AUDIO EQUIPMENT	85
FIGURE 5.15: H261 2.048MBIT/S N-ISDN.....	86
FIGURE 5.16: ISABEL MULTIMEDIA TERMINAL.	87
FIGURE 5.17: IP DATA TRAFFIC, PEAK 155.52 MBIT/S, MEAN 22 MBIT/S. (SDH TX LINE.)	92
FIGURE 5.18: VBR OPTIMISATION CONCEPT GRAPH	94
FIGURE 5.19: FTP OPTIMISATION EXPERIMENTAL SET-UP.....	95
FIGURE 5.20: MULTIMEDIA OPTIMISATION EXPERIMENTAL SET-UP	97
FIGURE 5.21: TIME TAKEN TO TRANSMIT A 5 MBYTE BINARY DATA BLOCK.	98
FIGURE 5.22: SHAPING AN FTP SESSION (DATA CHANNEL ONLY).	98
FIGURE 5.23: SHAPING AN ISABEL MULTIMEDIA TERMINAL	100
FIGURE 5.24: MAXIMUM SHAPER QUEUE LENGTH.....	101

FIGURE 5.25: SIMULATION SET-UP TO DETERMINE THE MULTIPLEXING DENSITY .	107
FIGURE 5.26: ALB FOR ISABEL MM TERMINAL COMPARED WITH ALB FOR CBR TRAFFIC	108
FIGURE 5.27: ALB FOR LARGER SWITCH BUFFER SIZES USING THE ISABEL MM TERMINAL.	109

Acknowledgements

The thesis has been conducted in parallel to work carried out on the European ACTS project EXPERT (AC094), but the work described in this thesis, except where indicated, is entirely the work of the author. I would like to thank the project and particularly WP4.1: Matthias Baumann, Stéphane Louis, Torsten Müller, Wouter Ooghe, Alfonso Santos and Martin Zeller for all there inspiring discussions and assistance on this topic. Thanks are also due to the ASPA staff in Basel, Switzerland, particularly Heinrich Blaser, Max Seiler, Georges Grun, and Urs Schenk for the use of the ATM testbed and their support.

I would like to thank my supervisor Prof. Laurie Cuthbert and Dr. John Schormans for all their hard work and assistance in making the PhD possible. Finally, I would like to thank my parents for their encouragement, the members of QMW's Telecommunication Laboratory for help and entertainment, Paul Seaman for his support and Tom Austin for the initial push.

Glossary

AAL	ATM Adaptation Layer
ABR	Available Bit Rate
ABT	ATM Block Transfer
ABT/DT	ABT/ Delayed Transfer
ABT/IT	ABT/ Immediate Transfer
ACR	Available Cell Rate
ALB	Admissible Load Boundary
ATC	ATM Transfer Capability
ATM	Asynchronous Transfer Mode
ATMF	ATM Forum
BCR	Block Cell Rate
B-ISDN	Broadband- ISDN
BT	Burst Tolerance
CAC	Connection Admission Control
CBDS	Connectionless Broadband Data Service
CBR	Constant Bit Rate
CD	Compact Disc
CDV	Cell Delay Variation
CDVT	CDV Tolerance
CER	Cell Error Ratio
CLP	Cell Loss Priority
CLR	Cell Loss Ratio
CMR	Cell Mis-insertion Ratio
CPN	Customer Premise Network
CRC	Cyclic Redundancy Check
CSCW	Computer Supported Co-operative Working
CTD	Cell Transfer Delay
DBR	Deterministic Bit Rate
DQDB	Distributed Queue Dual Bus
DTE	Data Terminal Equipment
DXI	Data eXchange Interface
FEC	Forward Error Check
FIFO	First In First Out
FTP	File Transfer Protocol
GoP	Group of Pictures
HDTV	High Definition Television
HEC	Header Error Check
HQ	High Quality
IAT	Inter-Arrival Time
IIP	International Inter-operator Portion
IP	Internet Protocol
ITP	International Transit Portion
ITU-T	International Telecommunications Union - Telecommunication
ISDN	Integrated Services Data Network
ISO	International Standards Institute

LAN	Local Area Network
LEX	Local EXchange
MBS	Maximum Burst Size
MCR	Minimum Cell Rate
MID	Multiplexing IDentifier
MM	MultiMedia
MMI	Man Machine Interface
MP	Measurement Point
MPEG	Motion Picture Experts Group
MTU	Maximum Transmission Unit
NIC	Network Interface Card
N-ISDN	Narrowband- ISDN
NNI	Network-Network Interface
NP	National Portion
NPC	Network Parameter Control
nrt-VBR	non-real-time VBR
OSI	Open Systems Interface
PCR	Peak Cell Rate
p.d.f.	probability density function
PDH	Plesiochronous Digital Hierarchy
PDU	Protocol Data Unit
PNNI	Public NNI
pt-to-ptCDV	Point-to-Point CDV
PVC	Permanent Virtual Circuit
QoS	Quality of Service
RM	Resource Management
RSVP	ReSource reserVation Protocol
rt-VBR	real-time VBR
SBR	Statistical Bit Rate
SCR	Sustainable Cell Rate
SDH	Synchronous Digital Hierarchy
SDU	Segment Data Unit
SECBR	Severely Errored Cell Block Ratio
SMDS	Switched Multimegabit Data Service
SN	Sequence Number
SVC	Switched Virtual Circuit
TAXI	TA eXchange Interface
TCP	Transport Control Protocol
UBR	Unspecified Bit Rate
UDP	User Datagram Protocol
UNI	User Network Interface
UPC	Usage Parameter Control
VBR	Variable Bit Rate
VCC	Virtual Channel Connection
VCI	Virtual Channel Identifier
VPC	Virtual Path Connection
VPI	Virtual Path Identifier
WAN	Wide Area Network

WWW World Wide Web

1. Introduction

1.1 Overview

ATM networks have been designed to integrate all of a user's communications needs including telephony, video distribution, LAN interconnect and IP services. [CUTH93], [ONVU95] and [PITT96] give a good introduction to ATM and Traffic Engineering. ATM uses a single physical medium and statistically shares resources in each direction of a link. Switching units interconnect links transporting cells between input and output; the units also provide a method to resolve contention onto the shared medium.

An application transmits traffic onto the network and this traffic is randomly mixed with other connections on the link. In order to provide guaranteed QoS some means of restricting the number of users on the network must be provided as a method to control the load: this is done by CAC.

The applications that share the resources are very different. Indeed, if the resources were shared equally amongst users, then applications requiring stringent bounds would be unable to function in the presence of bursty sources such as LAN interconnectivity. Applications need to be assessed so that the underlying network can provision for the particular service required. For example, a Circuit Emulation IWU would require a low-latency connection with a small variation in delay, while a data connection requires a loss-less medium. In answer to this, traffic categories with associated traffic differentiation methods have been defined in [ATM056]. Categories such as CBR, rt-VBR, nrt-VBR, ABR, ABT and UBR exist with the following methods of traffic differentiation: space priority, time priority, feedback and frame discard.

Undoubtedly, in the future, network users will demand greater and greater quality of service, [BELL96]. Users are exposed to the performance of an ATM network through the applications they use. Thus, users will grade the quality of service through subjective factors such as a good picture, clicks in sound and application freezing. This means that the quality of service is the measure of how the system,

(from the physical layer to application layer) performs as a whole, [MUST96]. Network performance parameters are used in quantitative methods to specify the required connection performance and, to allow a method to check whether the performance delivered is commensurate with the performance offered. However, users cannot determine these parameters directly from their applications; instead, they have to resort to test equipment to obtain measurements that provide CTD, CDV and CLR bounds.

Network performance and Quality of Service might well seem to be very different entities: one quantitative and one qualitative, but they are in fact related. Most users will be unable to assess the resources required or the network performance without expensive test equipment. Indeed, they will have to use subjective, “quality of service” to describe the application’s operation and map this to “best guess” network performance parameters.

The research presented within this thesis examines how quality of service can be provisioned over an ATM network. A discussion is presented that describes the ATM protocol stack in terms of services, connection objectives, and ATM transfer capabilities.

There are two main standardising bodies in the ATM arena; these are the ITU-T and the ATM Forum. Each has its own interpretation of Quality of Service and Network Performance, but comparisons can be made. [MUST96] describes a method to inter-work between the ITU-T and ATM Forum QoS parameters. This interworking is expanded within the background chapters of this thesis. Once the different layers have been discussed, mapping and choices will be made between the service classes, AAL layers and ATM transfer capabilities. Quality of Service very much depends on the selection of the protocol stack for a particular application, but an interface must be defined which allows different levels of service to be offered. [JUNG96] believes that the QoS of the application layer can be mapped into AAL layer PDU performance objectives; these parameters are further mapped into ATM network performance parameters. This is very difficult to achieve, as users would have to be quite technical to understand this mapping strategy. In addition, [DASI97] describes the differences in the performance of

identical end-systems. The fact that there is a substantial difference between end-systems would need to be accounted for when making the inter-layer mapping, particularly because this is a manufacturer dependent property. A single point of division is necessary to separate the responsibility of the user and the network operator: this is the ATM interface at the UNI. When different levels of service are offered, a measurement method must exist to determine if the guarantees specified have been met. These guarantees have been defined as network performance guarantees with six different performance parameters. A maximum of three of these can appear in a traffic contract negotiated at connection set-up, (Cell Transfer Delay, Cell Delay Variation and Cell Loss Ratio). Three others are used by the network operator to determine the accuracy of the network, (Cell Error Rate, Severely Errored Cell Block Ratio and Cell Mis-insertion Ratio). Measurement points, methods of calculating, and degradation of the different parameters are presented in this thesis.

Four ATM transfer capabilities have been defined with network performance guarantees: these are CBR, rt-VBR, nrt-VBR and ABR. The research in this thesis has covered CBR and rt-VBR ATM transfer capabilities, as these have the most stringent performance requirements. CTD, CDV and CLR have to be strictly defined for those services that require a high quality network layer. nrt-VBR and ABR have only CLR guarantees and are very dependent on the real-time services being transported over the network.

Real-time circuits are designed to carry the most stringent services. Hence, these circuits must have low latency, low cell delay variation and low cell loss. This thesis describes experimental research on switching devices and ATM networks to determine the maximum performance of the underlying network layer and to bound delays, variations, and losses caused by network elements. Mean queue sizes are estimated and the cell delay variation caused by queueing is examined. Delays are researched over local, international and intercontinental ATM networks. The networks researched are compared to standard methods for calculating performance measurements and concrete conclusions are drawn.

Once the underlying network has been bounded, services can be introduced. Network performance is compared to the applications' performance requirements to determine if an ATM network can accept the stringent services. In addition, an assessment of the applications' resilience to performance degradation is made. CBR services have the simplest traffic characteristics, but demand the greatest performance from the network. CBR applications tend to be based on existing telephony services, which need accurate synchronisation between end systems. Previous research has placed much emphasis on the buffering strategy of ATM switches. For instance, [LAND94], [KAWA96] and [LEE96] are a few of the authors that present different methods of reducing CDV and Jitter. These techniques add complexity to switching elements that need to operate at very high speeds. Moving the hardware complexity into the lower performance end-systems presents the cheapest solution for jitter rectification. Play-out buffers are present in all types of CBR equipment, and these are designed to absorb the CDV caused by ATM networks. The size of the play-out buffer is an important feature of the performance of CBR applications and, hence, QoS. This research examines CBR traffic and determines the optimum "play-out" queue size for a particular connection's bandwidth. Synchronous CBR traffic require high performance circuits. Provisioning for CBR ATCs means the most stringent services can be accommodated within an ATM framework. To guarantee QoS for all traffic, every real-time connection must be placed over these high quality circuits.

ATM networks have been designed bottom-up from the physical layer to the application layer. This means that while promoting a generic communication network, no particular application has been promoted. Thus, many other types of network can provide better and cheaper services for specific applications, [BELL96]. ATM needs to find a market niche that will allow it to be distinguished from other communication protocols. Certainly, the Internet Protocol will provide the "bread-and-butter" service over non-real-time circuits and the IP over ATM model will be prominent throughout this thesis. However, the greatest reason to provide ATM is the possibility to configure real-time services across a data network. The applications that are likely to use this real-time service category are video-on-demand services, interactive multimedia, and

multiparty games. Rt-VBR circuits pose an additional problem to the user and network operator, as preliminary research has shown that carrying such services makes inefficient use of network resources. As a consequence, the multiplexing density is determined within the research for a video and audio multimedia terminal. The traffic characteristics of such an application are altered with the introduction of a shaping device, and the multiplexing density is reviewed. Benefits and losses are reported when a shaping device is introduced to control, or change, the traffic profile.

A new QoS framework is proposed here that is simple to use, beneficial to customers and operators, and flexible enough to facilitate both synchronous and IP services. Careful consideration must be given to the three primary partners involved in service provision: the user, the network operator, and the service provider. Each will make gains from the communications exercise, and the service rendered must deliver a QoS better than (or at least as good as) existing networks to create the need for ATM. These gains can be maximised, and that is why a particular emphasis is placed on rt-VBR sources. The traffic characteristics can be modified using traffic shaping to provide cheaper services for the user and an increased number of connections of high-quality service categories. [GRAF97] describes a method of shaping based on multiple leaky buckets to reduce a video server's resource requirements. This argument is developed in the thesis into a method to efficiently transmit real-time sources to increase the benefits for the user and operator, while maintaining good end-to-end Quality of Service. By determining the responsible parties for network performance and QoS, a strategy can be adopted that will provide QoS to applications in a simple, efficient and cost effective manner.

1.2 The ATM Testbed

The ATM testbed used within the experiments was located in a Swisscom building in Basel, Switzerland and had originally been developed by RACE project R1022. Later, enhancements were completed by project R2061 and AC094. The testbed has been used for demonstrating new services, interconnection possibilities, and it has been used as the basis for many traffic

experiments designed to validate the considerable amount of theoretical work that has been undertaken on ATM traffic engineering. The ATM traffic experiments performed on the testbed have been divided into many areas, the principle ones being: Network Performance, Quality of Service, Traffic Management and Traffic Control. The testbed had access to one of the most modern telecommunications infrastructures in Europe. This European ATM network has been used successfully in many traffic experiments and demonstrations over European and global network scales.

1.3 Summary of Contribution

The research described in this thesis covers the area of Quality of Service in ATM networks. The emphasis has been placed on a workable QoS framework that can be provided in a public network environment. This work had been carried out within the scope of the RACE and ACTS projects of which the author was a member. These projects were:

- CEC RACE 2061 “EXPLOIT”
- CEC ACTS 094 “EXPERT”

These projects provided hardware facilities for the work, help and guidance through discussion, assistance in experimental data gathering, and funding enabling the work to be completed. However, the new concepts and ideas, and the experimental design were entirely developed by the author within the scope of the projects objectives. In addition, the author makes contributions that are beyond the scope of the project and are presented herein. The new contribution of the hypothesis is:

- A QoS mechanism that is able to guarantee network performance for all services. In particular, synchronous CBR services that require a stringent timing relation between end systems.
- As traffic characteristics are difficult to determine without test equipment, shaping devices can be used to bound the traffic parameters. Therefore, a

methodology for a QoS framework that provides a simplification of connection acceptance is proposed.

- Maintaining high performance levels on an ATM network, and allowing the user to mould the traffic with end system shaping enables user defined QoS levels. This QoS methodology enables the user the possibility to select an optimal QoS/Price function for the user's particular needs. In addition, in a multi-user environment each user is able to select their own independent QoS/price function.

1.4 Thesis Layout

In chapter 2, the ITU-T and the ATM Forum's standards are examined and a summary of the ATM protocol stack and the ATM transfer capabilities from a Quality of Service aspect is described. The Application, AAL and ATM layers are discussed with a particular emphasis on the selection of different AAL and ATC combinations to provide suitable QoS for applications. Chapter 3 examines only the ATM layer and describes the parameters used to define and measure network performance. It is these parameters that are used to guarantee network performance, and thus represent the contractual demarcation point. Chapter 4 describes in detail the objective of the hypothesis. It highlights the main principles for provisioning QoS, examines briefly the "community" of suppliers and users, describes the problem and offers a solution. Chapter 5 lists experiments and simulations used to determine the validity of the hypothesis. Experiments in Chapter 5 are configured on the ATM testbed, described in section 1.2. The simulations carried out in Chapter 5 use either the YATS [Baum96a] general purpose network simulator, developed in the EXPERT project, or special models were constructed by the author using the Simula programming language. Chapter 6 discusses the results found in Chapter 5 and Chapter 7 highlights the conclusion of the work.

2. The ATM Protocol Stack

2.1 The Application Layer

Currently there are four specified QoS classes (and one unspecified), which are used to provide different services over ATM. These have been defined by the ATM Forum. An example selection of applications has been given in table 2.1 to give a perspective of the type of applications likely to be used within each service class. This is not an exhaustive list as the application category depends also on the traffic characteristics. Thus, for example, High Quality Audio could be placed in either the CBR or rt-VBR sections depending on how the application has been constructed.

Class 1 Circuit Emulation CBR Video	Class 2 VBR Audio and Video	Class 3 Connection Orientated Data Transfer	Class 4 Connectionless Data Transfer
N-ISDN Circuits	Packetised Video	Pager Messaging	Datagrams
Emergency Service Applications	Audio Retrieval	Voice Box Messaging	File Server Access
Online Video	Teleshopping	Teletext	Email (IP)
Online Audio	Pt to Pt MM Terminals	Minitel	Net News (IP)
HQ Audio	Multiparty MM Terminals	LAN Emulation, Frame Relay, X.25	Telnet (IP)
Videophones	Online Video Shops		WWW (IP)
Pt-to-Pt Telephone	Online Computer Games		FTP (IP)
Multiparty Telephone	Long Distant Feedback Control Loops		SMDS
Mobile Telephone	MPEG 2		CBDS
Mobile Data			
Television Broadcast			
HDTV			

Table 2.1: Application with QoS Classes

All applications are different, requiring different bit rates, cell loss ratios and cell transfer delays for example. Therefore, the traffic characteristics are very

dependent on the way in which an application has been designed and developed. Each application displays different characteristics. Theoretically, this means that each application could have its own particular set of UPC parameters to use on connection set-up and its own QoS requirements, but this is not practically possible as the calculation of the overall network performance and the method of charging would be too complex for the network operator and customers alike.

2.2 The AAL Layer

The main standardisation organisations: ITU-T and ATM-Forum differ in their definition of the QoS Classes. Both notions are complementary and relationships between them can be established.

The ITU-T approach for the definition of Quality of Service classes can be interpreted as being driven by the requirement of the ATM layer; a bottom up methodology. In this approach, the quality of the service offered at the ATM level is defined. These appear at the bottom of figure 2.1 and are numbered from 1 to “U” (U for unspecified) meaning, schematically, from a well-specified, and hence guaranteed, Quality of Service, down to a non-specified Quality of Service with poor performance.

The ATM-Forum Quality of Service classes are specified in terms of sets of application types. They are numbered 1 to 4 at the top of figure 2.1. This numbering does not specify an increasing or decreasing level of quality, as each class addresses different applications types. This top-down approach is the opposite of the bottom-up one of the ITU-T. Therefore, the ATM Forum methodology can be interpreted as being driven by the requirements of the applications.

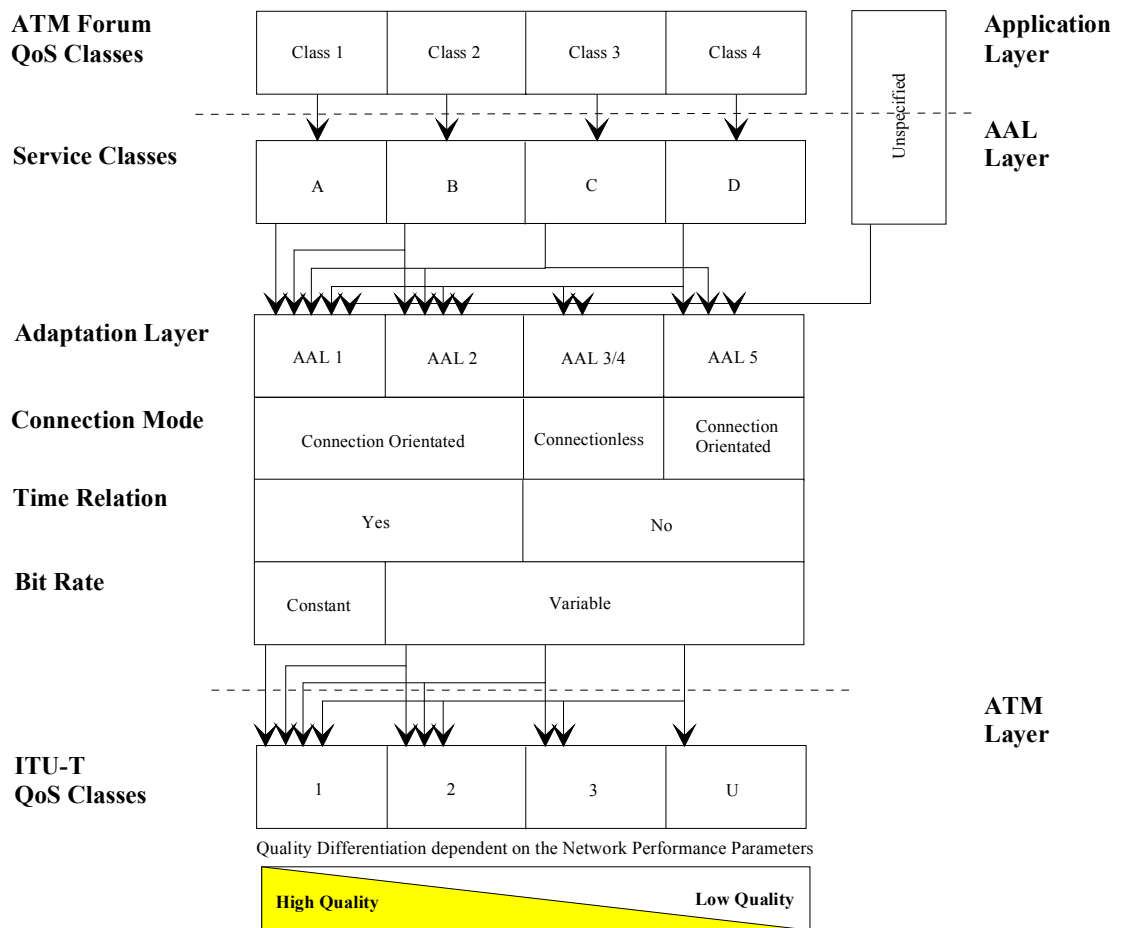


Figure 2.1: The ATM protocol stack and QoS Classes

Figure 2.1 represents the relationships that can be established between these two representations of the QoS classes.

The choice of an AAL for a certain application is determined by how the AAL characteristics fulfil the application needs or requirements. These are mainly classified into three criteria:

- Existence of a *time relationship* between the source and the destination. Either there is a time relationship need, for example, with voice applications, or there is no timing relation, as for file transfer.
- Characteristics of the *bit rate* associated with the data transfer: constant (as for NISDN circuit emulation) or variable (as are many multimedia terminals).

- The *connection mode* is either connection-oriented (as are any applications using video and audio media) or connectionless, (for example, Connectionless Broadband Data Service or Switched Multimegabit Data Service).

Hence, the AAL Service Classes issued from the combination of these criteria, can be associated with application types:

- Service Class A suits circuit emulation and constant bit rate video and audio e.g. telephony, MPEG images retrieval
- Service Class B suits variable bit rate video and audio e.g. variable bit rate MPEG 2
- Service Class C suits connection-oriented data transfer e.g. X.25, Frame Relay
- Service Class D suits connectionless data transfer e.g. CBDS/SMDS

These are defined in [I.362]. Their role is to enhance the service by the ATM layer to the higher layers. The functions performed within the AAL depend on the requirements of the higher layers.

The complete description of the AAL can be found in [I.363].

2.2.1 ATM Adaptation Layer 1

The service provided by AAL1 to the higher layers consists of transferring data units with bit rates kept constant during the lifetime of a connection. It transfers timing information between the source and the destination. AAL1 provides a means to indicate unrecoverable losses and errors to higher layers.

Utilising the Sequence Number (SN), and an additional parity check cell, interleaving, which provides an effective method of forward error correction (FEC) can be implemented. However, the introduction of this method of FEC increases the packetisation delay.

2.2.2 ATM Adaptation Layer 2

The service provided by AAL2 to the higher layers encompasses the transfer of data units at variable bit rates. It also transfers timing information between the source and the destination. AAL2 provides the means to indicate unrecoverable losses and errors to higher layers.

Recent standardisation work has defined AAL2 procedures. Previously, it was thought that AAL2 would be similar to AAL1, except that it would be transferred through the rt-VBR ATC. The new format for AAL2 includes a channel identifier to enable the multiplexing of AAL frames (16 to 64 bytes) into an ATM stream identified by a single VCC or VPC. Thus, larger AAL2 frames will be allowed to break over ATM cell boundaries or several small frames can be carried in one cell, unlike AAL3/4 PDU. The main use of this type of service would be to carry low bandwidth mobile telephony type services where packetisation delay is a problem. A full description can be found in [I.363].

2.2.3 ATM Adaptation Layer 3/4

The service provided by AAL3/4 consists of transferring variable length data units from one source to possibly several destinations via connectionless servers. Originally, AAL3 was connection-oriented and AAL4 connectionless. They are now merged and the service provided is connectionless.

The main characteristic of AAL3/4 service is that it works according to two modes: streaming mode or message mode. The difference between the two is only a question of reducing packetisation delay, which should not influence intermediary switches or servers. Streaming mode sends service data units received from the ATM layer to the upper layer as soon as they arrive. Message mode waits for all the service data units belonging to an upper layer message to arrive and re-assembles the PDU, before passing the message as a single data unit to the upper layer.

In the case of the message mode, the service is offered at the expense of an additional delay (because of the re-assembly of the data unit into a message). In

the case of the stream mode, the service is offered at the expense of the message validity before passing it to the upper layers.

Using the multiplexing identification (MID) function of the transported messages, AAL3/4 can multiplex information onto a single ATM connection. AAL3/4 will be used exclusively for data transfer.

2.2.4 ATM Adaptation Layer 5

The specification of AAL5 was a response to the criticism of AAL3/4. The latter (AAL3/4) is intended to transport data in a connectionless fashion. It requires, however, the utilisation of a significant part of the ATM cell payload: 4 out of 48 bytes. AAL5 is the simplest AAL service. The service information necessary for the recognition of transported information is found only at the tail-end of the data unit (in the last ATM cell transporting the segmented message): the PTI bit in the ATM header indicates whether the information transported in the payload of the cell is the last part of the message to be passed to the higher layer.

The service provided by AAL5 is very similar to the service provided by AAL3/4, except that the AAL3/4 method of multiplexing cannot be supported. Therefore, broadcast channels cannot be implemented within AAL5. AAL5 particularly suits data transfer in the context of a heterogeneous network inter-connection (i.e. Frame Relay over ATM)

2.3 The ATM Layer

Six categories have been defined within the ATM layer. These service categories relate traffic characteristics and QoS to network performance guarantees. There are three real-time service categories CBR, rt-VBR and ABT¹. The traffic descriptors contain only the Peak Cell Rate (PCR) or both the PCR and the Sustainable Cell Rate (SCR). Three non-real-time categories exist: nrt-VBR, ABR and UBR. Each ATM transfer capability (ATC) can be applied to both

¹ ABT has an additional Block Cell Rate parameter that specifies the PCR for a particular block transfer.

VPCs and VCCs. Resources management cells (RM-cells) are used to control an ABR connection; these cells are allowed within CBR, rt-VBR, nrt-VBR and UBR. If RM-cells are present in these ATM categories then the RM-cells are considered a part of the user's traffic stream.

2.3.1 Constant Bit Rate (CBR)

The CBR service category is used for connections that require a fixed amount of bandwidth continuously available during the lifetime of a connection. CBR is intended to support applications with a stringent time relationship and bounded time delay and time delay variations. This is representative of applications such as telephony and more generally Circuit Emulation.

The QoS is maintained at the same level for all the cells transmitted during the lifetime of the connection, while the cells conform to the traffic contract. The Peak Cell Rate and the Cell Delay Variation Tolerance define the traffic contract. PCR represents the highest rate at which cells can be sent during the connection lifetime and CDVT is the upper limit of the allowable cell delay variation. The CBR source may transmit the cells at any rate up to the PCR and may even be silent for periods of time. Cells that have been delayed a greater length of time than the maximum delay specified (maxCTD) have a significantly reduced value for the application.

As indicated in table 2.2, this Service Category belongs to the ATM-Forum QoS Class 1. Applications generating traffic of this type are likely to be implemented over AAL1. The ITU-T QoS Class 1 is likely to be ordered.

2.3.2 Real-time Variable Bit Rate (rt-VBR)

The service category rt-VBR is intended to describe sources that transmit at a rate that varies in time. These are real-time applications, with stringent speed constraints, i.e. delay and delay variation. Applications in the rt-VBR category are based on voice and video, for example, video conferencing. This type of application does not support a delay higher than 25 ms because of the audio component. The real-time constraint should guarantee a synchronisation of voice

and image and guarantee an accurate representation of the transmitted components.

The rt-VBR service category is characterised in terms of Peak Cell Rate, Sustainable Cell Rate and Maximum Burst Size (MBS) allocations. These values define the boundaries of the traffic characteristics. Cells that do not conform to these requirements are assumed to be of a significantly reduced value to the application. These cells may be discarded or simply tagged. In both cases, actions must be taken at each level of the protocol stack in order to recover from this state. The cell rate is expected to vary with time and the source can be described as “bursty”. Real-time VBR may support statistical multiplexing of real-time sources.

Applications in this service category are likely to be implemented using AAL2. However, AAL2 is far from being deployed on a large scale and it is likely to be redefined for mobile services. AAL1 or AAL5 could be considered suitable substitutes.

2.3.3 Non-real-time Variable Bit Rate (nrt-VBR)

The definition of the traffic parameters attached to the non-real-time VBR service category is identical to those of the rt-VBR. The applications using this service category expect lower cell loss ratios, but do not have any speed definitions.

The nrt-VBR service category appears best suited for statistical multiplexing gain. The indeterminate time constraints give the possibility of using a large buffer along the connection in order to achieve this goal. Due to the absence of time constraints, it can be foreseen that the implementation of this service category could use AAL3/4 or more probably AAL5.

2.3.4 Available Bit Rate (ABR)

The Available Bit Rate service category is characterised by the existence of a mechanism allowing the network to vary the transfer characteristics from the connection set-up throughout the duration of the call. This implies that the

network resources particular to the connection may not remain constant during the connection lifetime.

The variations managed by the traffic control mechanisms are reported to the source application by the means of feedback traffic, which aims to control the source rate. Compliance to the variations from the feedback signal should guarantee a low CLR for the application and a fair share of the available bandwidth.

There is no time constraint defined in terms of delay or delay variation boundaries, and ABR service categories are not intended to support real-time applications. PCR and MCR allocations set the boundaries within which the resources provided by the network have to remain. The MCR may be specified as zero, in which case, the connection could stop until network resources become available. The bandwidth available from the network, defined as the ACR can vary between the PCR and MCR. The source will always be able to transmit at the MCR despite possible contradicting information from the RM-cells.

It is necessary to have large buffers to offer an ABR service over a network; this is due to the bursty nature of the traffic types and the delay, caused by network latency, in RM-cells that reduce a source's bandwidth. It is likely that the ABR service category will make use of the AAL5.

2.3.5 Unspecified Bit Rate (UBR)

The UBR service category is intended for non-real time applications, i.e. without stringent delay and delay variation constraints. This applies to computer networks in general and file transfer (FTP for instance), E-mail or Telnet sessions.

The characteristics of the UBR service category determine that no commitments are given with respect to the QoS objectives defined in the other service categories. A network may, or may not, apply PCR to the CAC and UPC functions.

PCR allocation is still possible without any guarantee of conformance parameters from the network. Networks that do not apply PCR use this value for information purposes only. Loss and error recovery or congestion control mechanisms are implemented in the application, and not at lower network layers. It is probable that the UBR service category will offer best effort services only.

2.3.6 ATM Block Transfer (ABT)

ABT capability can be interpreted as a “non-permanent” CBR service category. For each block of information sent, a resource management cell reserves the bandwidth required. The bandwidth and resource allocation are performed for a single block only.

During a block transfer, the traffic characteristics and QoS requirements are the same as during the lifetime of a connection utilising the CBR service category. The QoS experienced by the cells within a block is comparable, then, to the QoS experienced by the cells of a CBR connection.

The connection is not released after each block transfer. The network has, therefore, to allocate sufficient resources to ensure that the QoS is kept equivalent to all the cells within the block. The bandwidth is allocated based on a Peak Cell Rate allocation. The cell rate ordered for a block transfer is called a Block Cell Rate. The negotiated element at the connection establishment is then the “absolute” maximum PCR and CDVT and the frequency of the re-negotiation of the BCR, (“Absolute” meaning that the negotiated PCR is the upper bound for the intermediary ordered BCR).

Two subsets of ABT are: ABT with immediate transfer (ABT/IT) and ABT with delayed transfer (ABT/DT). In the first case, the blocks are transferred towards the network before receiving an acknowledgement regarding the availability of the required BCR. This may be at the expense of block losses. In the second case, the block is sent only after the BCR has been confirmed, i.e. agreed between the user and the network. This is done at the expense of the transfer delay.

ABT/IT is likely to use AAL3/4 while ABT/DT is likely to use AAL5.

2.4 Traffic Control Functions

ATM networks are designed to support a variety of services. Quality of Service cannot be provided to ATM connections without some form of traffic control. The network must provide QoS to network applications and this is achieved by providing appropriate means to differentiate between traffic types. The main objective of traffic control is to protect the network and the end systems from congestion in order to maintain network performance guarantees. In addition, traffic control is used to promote the efficient use of network resources.

2.4.1 Connection Admission Control

Signalling in the call set-up phase is used to indicate the amount of resources required by the application. Connection Admission Control (CAC) is defined as the set of action taken by the network to assess whether the connection should be admitted or rejected. CAC is also used in connection re-negotiation to determine if the additional resources required could be allocated to the user or not.

2.4.2 Feedback

The network and user regulate the amount of traffic onto the network with feedback. Signals are sent to the end terminals from the network to alter the amount of traffic on to the network. This signal is varied according to network conditions.

2.4.3 Policing

Policing is defined at two locations; Usage Parameter Control (UPC) is traffic control at the User-Network Interface (UNI) and Network Parameter Control (NPC) is traffic control at the Network-Network Interface (NNI). The network monitors the traffic parameters, negotiated at the call-set-up of a single ATM connection, with a policing function. This connection could be either a VPC or a VCC. The main purpose of this traffic control function is to protect the network from malicious as well as unintentional misbehaviour that can affect the QoS of other already-established connections. Once unacceptable conditions are detected

the policing function can either discard or tag the offending cells within the connection.

2.4.4 Priority

Space priority is useful within some service categories. The end user can set a cell loss priority (CLP) in the traffic stream for cells that have a less importance to the user. The network may treat these markings as transparent or as significant.

When the markings are significant, the network can selectively discard cells with low priority upon switch congestion. This helps maintain the QoS of the remaining connections.

Time priority is an important mechanism to increase the utilisation of an ATM network. ATM is an integrated services digital network, and data services are naturally very bursty. Thus, to increase the utilisation, very large buffers are required to absorb the bursts. Large buffers reduce, however, the ability of the network to provide time critical services. Therefore, a priority mechanism based on separating traffic into different queues is necessary to maintain a differentiation between time critical and non-time critical applications.

2.4.5 Network Resource Management

The service architecture allows logical separation of connections according to service characteristics. Cell scheduling and resource provisioning are switch and network specific. These can be used to provide isolation and access to resources. VP networks are a useful tool for managing resources.

2.4.6 Frame or Packet Discard

A network will always lose cells, particularly if the cells have a low priority. Many types of data networks transmit the information in the form of frames. When a cell belonging to a frame is discarded, the data frame is rendered invalid. Using a frame-discard policy, a cell loss will cause the removal of all of the cells from that particular frame. This allows resources to be released for other

connections. When a whole frame is discarded, the responsibility of recovery is left to high layer protocols.

2.4.7 ABR Flow Control

ABR flow control allows the bandwidth to be allocated dynamically. The bandwidth is shared amongst the participating users in a fair manner.

2.5 ITU-T and ATMF QoS Mapping.

Quality of service has a different meaning according to the context in which it is defined. It depicts the accuracy of a service offered to a user. The user can be a person or a protocol layer. The service can be a high layer service like video application or a lower layer service (as defined by OSI). Hence, the Quality of Service can only be defined within a specified context.

In the OSI architecture, QoS is specifically defined for the transport layer i.e. between the session layer and the network layer. As an intermediary layer, the transport layer “proposes” to the session layer five Quality of Service classes (four specified and one unspecified), depending on the reliability of the underlying network. QoS classes range from zero (when the underlying network is reliable), up to four (which contains the maximum control and recovery functions).

The transmission of data has become more reliable with the introduction of digital technology (rare losses rather than erratic error conditions for instance). The notion of Quality of Service has moved up the ISO protocol stack from the physical layer. Quality of Service is driven by both the requirements of an application and the guaranteed network performance provisioned by the network provider. Thus, the errors have become a subject of correct dimensioning of control parameters over the network. The reliability of a network relies on good management of resources, guided by appropriate dimensioning of the network control parameters.

The dimensioning of these network control parameters requires an accurate study of the network performance parameters themselves. Some essential network

performance parameters need to be gathered into sets or Quality of Service classes. Within each class, objectives for these performance parameters are defined. These objectives define the boundaries of the Quality of Service.

2.6 Service Categories and Quality of Service Classes

As stated in [ATM056], the ATM-Forum QoS classes presented in table 2.2 can be matched with the AAL service classes. There are four specified QoS service classes that have a one-to-one match with the AAL service classes. An additional so-called “unspecified” QoS class is defined.

ATM-Forum QoS Class	Application Nature
Class 1	Requiring performance comparable to current digital private line performance
Class 2	“Packetised” video and audio in teleconferencing and multimedia applications
Class 3	Inter-operation of connection oriented protocols such as Frame Relay
Class 4	Inter-operation of connectionless protocols, such as IP or SMDS
Unspecified	The support of the “best effort” services

Table 2.2: QoS Classes mapped to applications, as recommended by the ATMF

Besides the unspecified QoS class, the definition of these classes is tied to objective settings regarding the network performance parameters. The Quality of Service is defined in terms of the applications that they encompass. Moreover, the network performance objectives are defined according to the needs of envisaged applications.

A QoS Class, as defined by the ITU-T in [I.356], reflects the quality of the service at the ATM level. The objectives of a set of specific network-performance parameters determine the accuracy of the service offered from the ATM layer to

the AAL. The definition of these QoS classes is, however, made in total independence of the type of “bit rate” produced by the combined application-AAL implementation on top of the ATM layer. Nevertheless, ITU-T in [I.356] recommends QoS Classes be chosen according to the nature of this bit rate. These recommendations are presented in table 2.3.

ITU-T QoS Class	Bit Rate Type
class 1 stringent class	DBR
class 2 tolerant class	SBR, ABR, ABT
class 3 bi level class	SBR, ABR
U class ²	UBR or any other bit rate nature

Table 2.3: QoS Classes mapped to ATCs as recommended by the ITU-T

The QoS classes are, as defined in terms of network performance parameter objectives, presented further on in this chapter.

From figure 2.1 and tables 2.2 and 2.3, it is possible to make some recommendations. Indeed, for a given ATM-Forum QoS Class (1 to 4, U) to which an application belongs, the choice of the ITU-T QoS Class is either unique or multiple. The choice of protocol stack functionality is something that cannot be imposed, but only recommended. Besides the performance, requirements and characteristics, economics may impose the greatest constraint.

The various bit rate types detailed in table 2.4 correspond to the ATM-Forum definitions of Service Categories. These are defined in [ATM056]. The ITU-T defines these various bit rates in terms of ATM Transfer Capabilities (ATC). ATC and Service Categories can roughly be matched for a large part. The naming conventions are however different. Table 2.4 presents the implicit correspondence of Transfer Capabilities with Service Categories.

² U stands for Unspecified

ITU-T ATM Transfer Capability	ATM-Forum Service Category
Deterministic Bit rate (DBR)	Constant Bit Rate (CBR)
-	Real-time Variable Bit Rate (rt-VBR)
Statistical Bit Rate (SBR)	Non-real-time Variable Bit Rate (nrt-VBR)
Unspecified Bit Rate (UBR)	Unspecified Bit Rate (UBR)
Available Bit Rate (ABR)	Available Bit Rate (ABR)
ATM Block Transfer (ABT)	-

Table 2.4: Mapping ATCs between the ITU-T and ATM-Forum

2.6.1 Service Categories

Both standardisation bodies, ITU-T in [I.371] and ATM-Forum in [ATM056], have identified different bit rate types at the ATM layer. ITU-T defines ATM Transfer Capabilities whilst the ATM-Forum defines Service Categories. Both organisations define quite similar traffic types. The ATM Forum divides the variable bit rate model in real time and non-real time variable bit rate. This illustrates that the approach of the ATM Forum is mainly driven by the applications' needs.

The Service Categories are defined by specifying

- Network performance parameters.
- Traffic contract parameters.

2.6.2 Definition of Service Categories with Network Performance Parameters

ATM layer network performance parameters are used to characterise the Quality of the Service at the ATM level. The parameters used in the traffic contract to scale the level of service supplied within service categories are listed below.

- *Cell Loss Ratio*, which is the ratio of lost cells to transmitted cells.

- *Maximum Cell Transfer Delay*, which is the maximum time for a cell to be transferred from its origin to its destination
- *Peak-to-peak Cell Delay Variation*, is the difference between the shortest CTD (best case and fixed value) and the longest CTD (worst case minus the probability of being lower than a certain threshold). This parameter holds the notion of the time relationship evoked previously in this chapter.

A simplified presentation of the Service Categories specification is shown in table 2.5. Beside the Quality of Service parameters mentioned above, this table also uses traffic parameters required to specify the Service Categories:

- Peak Cell Rate and Cell Delay Variation Tolerance.
- Sustainable Cell Rate and Maximum Burst Size.

These traffic parameters are the elements commonly referred to as belonging to the traffic contract negotiation.

ATM Attributes	ATM Layer Service Categories					
	CBR	rt-VBR	nrt-VBR	ABR	UBR	ABT
QoS Parameters						
CLR	specified	specified	specified	specified	unspecified	not applicable
peak-to-peak CDV	specified	specified	unspecified	unspecified	unspecified	not applicable
maxCTD	specified	specified	unspecified	unspecified	unspecified	not applicable
Traffic Parameters						
SCR, MBS, CDVT	not applicable	specified	specified	not applicable	not applicable	unspecified
PCR, CDVT	specified	specified	specified	specified	specified	specified
Other Parameters						
feedback ³	unspecified			specified	unspecified	

Table 2.5: ATM Attributes and Service Categories specifications

2.7 Application Categories, AALs and ATM Service Categories

2.7.1 Overview

The mapping between application categories and ATM adaptation layers (AAL), and the mapping between AALs and ATM service classes depends on three classes of constraints:

1. Services provided by the respective lower layer.

In particular, AALs provide certain functionality (e.g. establishment of timing relations or assured mode transmission) which have to fit with the application needs.

2. Quality of service commitments provided by the lower layer.

If, according to point 1, the functionality has been selected, the quality of the

³ This concern the part of the management traffic a user receives from the network (see section describing ABR).

services provided (in terms of speed, accuracy, and dependability) has to be assessed.

3. Pricing and limitations of traffic control functions.

Due to the implications of network utilisation, different ATM service categories will follow different pricing strategies. Additionally, parameter boundaries of traffic control functions will limit the possible choices.

In general, a complete framework of possible mapping scenarios cannot be given. New applications with specific combinations of transport demands and new higher protocol layers will appear. The specification of ATM service categories – especially for the needs of computer communication – is in progress, and the definition of service specific parts of the AAL is not yet complete.

2.7.2 Functionality Mapping Requirements

2.7.2.1 Selecting an AAL

The two most important services influencing the choice of an AAL are provision of timing relations between the source and destination, and the forms of error handling. According to ITU-T recommendation [I.362], the implementation of connectionless transfer functions will be on top of the AAL, so that the connection mode is of less importance to the AAL. This is reflected by the combination of separated AAL layers 3 and 4 into a single specification. Capabilities for structured data transfer (AAL1) and support of higher layer connection multiplex (AAL2 and AAL3/4) are other criteria.

Many applications relying on timing relations, and therefore use AAL1 (ITU-T recommendation [I.363]); these include circuit emulation, high quality video signal transport, and voice band signal transport. Such applications belong to the group already effectively supported by existing circuit-switched networks. Since they often have been developed under the assumption that a synchronous transport medium is used, they rely on the implicit timing relation provided by these networks. The emulation of this capability is therefore mandatory. This does not hold, however, for the group of real-time computer applications, such as MPEG

video transmission or computer-based multimedia sessions. Since software protocol stacks always introduce variable and unpredictable delays (depending, for example, on the machine load), these applications implement their own mechanisms to establish timing relations. Another factor is that compression techniques like MPEG maintain a receiver buffer in order to allow forward and backward references needed for the compression algorithm. This receiver buffer can probably be used to handle transmission delay variations. Consequently, computer applications with real-time requirements often can use AAL3/4 and 5, given that the transfer delay variations of the underlying ATM transfer capability remain appropriately bounded.

The second large group of computer applications comprises interactive text/data/image transfer, messaging and retrieval. They can accommodate virtually arbitrary transmission delay variations and, therefore, should use AAL3/4 or AAL5. The functionality of multiplexing several higher layer connections into one ATM layer connection is seldom needed. Thus, AAL3/4 with the relatively large overhead of four bytes per cell is not often applied in practice.

Two kinds of error handling are currently included in the AAL specification. Forward-error correction in combination with interleaving is defined as part of the service specific part of AAL1, and is foreseen for applications like high quality video signal transport or circuit emulation. Again, it is intended to emulate as much as possible the behaviour of circuit-switched transport networks. When AAL1 is mapped to CBR, and CBR traffic in turn is multiplexed together with rt-VBR, a slight overbooking of a network link could result in a very low but non-zero cell loss probability. To maintain source clock recovery, AAL1 in such cases would include invalid data (ones or zeros) into the traffic stream. Forward-error correction would be a way to reconstruct the payload of lost cells.

The assured mode of AAL3/4 and 5 will be realised by the service specific part of the AAL convergence sub-layer. Nevertheless, this part is not yet fully specified [I.363]. Furthermore, it remains questionable as to whether it is useful to complement retransmission protocols of the higher layers that act in practice on

the same PDUs. However, assured mode could be useful for future applications that want to use “pure” ATM with guaranteed error-free transmission.

2.7.2.2 Mapping to ATM Service Classes

To maintain a source clock recovery, both the cell delay variation and the cell loss ratio have to be bounded. Therefore, only CBR and rt-VBR can be used to transport traffic of type AAL1 and AAL2.

For all other AAL types, the mapping is only restricted by QoS, pricing, and traffic control issues. As already mentioned, this possibly may not hold in the case where real-time computer applications are mapped onto AAL3/4 or 5. Since neither, nrt-VBR, ABR, nor UBR gives any guarantees about transmission delays and their variations, they are inappropriate under these circumstances.

2.7.3 Quality of Service Mapping Requirements

The ITU-T recommendation [I.350] introduces a general framework for QoS assessment. A 3×3 matrix given in table 2.6 defines nine generic performance parameters.

Function	Speed	Accuracy	Dependability
Access	Access Speed	Access Accuracy	Access Dependability
Information Transfer	Information Transfer Speed	Information Transfer Accuracy	Information Transfer Dependability
Disengagement	Disengagement Speed	Disengagement Accuracy	Disengagement Dependability

Table 2.6: Generic performance parameters for QoS assessment

Speed is the performance criterion that describes the time interval that is used to perform the function or the rate at which the function is performed with the desired accuracy.

Accuracy is the performance criterion that describes the degree of correctness with which the function is performed with desired speed.

Dependability is the performance criterion that describes the degree of certainty with which the function is performed regardless of speed or accuracy, but within a given observation interval.

For our purposes, only the criteria related to information transfer will be considered. The association of detailed QoS parameters to these classes is different for each application (especially for subjective assessment parameters). The following classification of the ATM layer performance parameters has been given [ATM810].

QoS Performance Parameter	QoS Assessment Criteria
Cell Error Ratio	Information Transfer Accuracy
Severely-Errored Cell Block Ratio	Information Transfer Accuracy
Cell Mis-insertion Rate	Information Transfer Accuracy
Cell Loss Rate	Information Transfer Dependability
Cell Transfer Delay	Information Transfer Speed
Cell Delay Variation	Information Transfer Speed

Table 2.7: ATM layer QoS parameters and generic assessment criteria

The relationships between transfer speed, transfer accuracy and transfer dependability on different layers is many-sided. Suppose, for example, the ATM layer relies on a bit-error-free transmission system, i.e. accuracy and dependability are optimal. Then the cell transfer accuracy on the ATM layer also will be optimal. Statistical multiplexing decreases the dependability (cell loss) and influences cell delay (that is seen as a speed measure), both CDV and CLR can degrade the transfer accuracy (inclusion of ones or zeros as dummy bytes) and dependability (loss of synchronisation, [MERA96]) for AAL1 and subsequent audio transmission. For an FTP connection, cell loss (dependability) maps to AAL5 SDU loss due to CRC check errors (dependability) and eventually onto transfer speed or even dependability in the case of connection time-outs, while accuracy will stay excellent. For MPEG video transmission over AAL5 as another

example, AAL5 SDU loss will result in lower accuracy (image corruption) or even lower dependability (decoder crash).

2.7.3.1 Selecting an AAL

To decide whether the QoS provided by an AAL is sufficient, the mapping between AAL QoS and user perceived QoS has to be investigated. The most important AAL QoS measures are:

- Dependability of synchronisation and accuracy of transmission (bit error rate due to dummy byte inclusion) for AAL1.
- Frame loss ratio for AAL3/4 and 5.

Unfortunately, their mapping onto user perceived QoS depends on a variety of factors, such as higher layer protocols, compression techniques, real-time requirements and subjective human assessment. Higher layer protocols implement their own flow control mechanisms that may interact with ATM layer mechanisms. Issues, such as, selective / unselective retransmission, fast / slow start of transmission, and techniques to observe network status (congested / uncongested) influence significantly the behaviour of the combination of ATM traffic control and higher layer protocols. Compression techniques in general reduce the bit rate needed, but also make the traffic stream more vulnerable to loss or data corruption. The internal structure of the compressed data stream may lead to distinct results in error propagation dependent on the time instant of data loss or corruption. As an example of the problem of subjective quality assessment, compare image corruption during a scene with high movement and during a scene-change. Both will produce high bit rates if a VBR video coder is used. While the quality degradation could be less obvious during a high-movement scene, starting a new scene with a corrupted image would be disturbing as the error may propagate through succeeding picture frames.

2.7.3.2 Mapping to ATM Service Classes

ITU-T recommendation [I.365] introduces provisional definitions that relate the boundaries of the ATM layer QoS performance parameters to four QoS classes. Table 2.8 depicts the boundaries given for cell transfer delay (CTD), 2-point cell transfer delay variation (2-pt CDV) and cell loss ratios. CLR is for high and low priority cells (CLR_{0+1}), and for high priority only (CLR_0), cell error ratio (CER), cell mis-insertion ratio (CMR), and severely errored cell block ratio (SECBR).

	CTD	2-pt. CDV	CLR_{0+1}	CLR_0	CER	CMR	SECBR
Default objectives	no default	no default	no default	no default	$4 \cdot 10^{-6}$	1/day	10^{-4}
Class 1 (stringent class)	400 ms	3 ms	$3 \cdot 10^{-7}$	none	default	default	default
Class 2 (tolerant class)	unspec	unspec	10^{-5}	none	default	default	default
Class 3 (bi-level class)	unspec	unspec	unspec	10^{-5}	default	default	default
U class	unspec	unspec	unspec	unspec	unspec	unspec	unspec

Table 2.8: ATM layer QoS classes proposed by [I.365]

To decide which ATM service category to use for a certain AAL, the committed CLR and CTD/CDV values can be checked. Besides this, a very important issue is the mapping of SDU loss between ATM layer (CLR) and AAL (frame loss ratio, FLR), which is mainly determined by the distribution of cell losses in time, which in turn is influenced by traffic control functions of the ATM service classes. For CBR, which is not subject to statistical multiplexing, cell losses and bit errors will be distributed more or less uniformly in time. This leads to frame losses for every cell loss or bit error, because every failure causes an AAL CRC error in another frame. The opposite will be observed with techniques like early packet discard (under discussion in the UBR service category), which tries to increase the ratio of CLR and FLR by selective, frame dependent cell removal. Cell losses caused by statistical multiplexing (VBR) also tend to be strongly correlated

[BLON91]. The close correlation of lost cell again results in a better ratio of CLR and FLR than for CBR.

2.7.4 Mapping Pricing and Traffic Control

2.7.4.1 Pricing

Pricing will probably have one of the most important impacts on the selection of an ATM service category. From the network provider's point of view, both the achievable network utilisation and the effort to maintain the necessary traffic control functions have to be considered. Although, for example, the costs of control functions for ABR and ABT in wide area networks cannot yet be estimated, the following overall relations may be predicted:

1. CBR

The relatively highest price (related to the peak bit rate) of this service category results from two factors: the high quality channel provided is virtually transparent and fully available throughout the lifetime of the connection.

2. VBR

As with CBR, VBR provides a transparent channel with high quality. Relative to the peak bit rate, it will be less expensive since multiplexing gain can be achieved. This gain can stem from VBR multiplexing based on appropriate source policing and buffer or bandwidth reservation, or it may be achieved in combination with ABR/ABT/UBR.

3. ABR and ABT

These service categories allow "gaps" in bandwidth to be filled by VBR/CBR connections, when used in an integrated environment. If applied on private LANs, the link utilisation can be relatively high. Additionally, the commitment on available bandwidth is less stringent from the user's point of view. Therefore, a lower price - compared to VBR/CBR - should be possible. However, a considerable amount of traffic control will be necessary to operate ABR/ABT in wide area networks, so increasing the cost of terminal equipment.

4. UBR

Since no QoS commitments are made, this service category should be the cheapest. Additionally, the open-loop traffic control strategy allows for yet higher link utilisation than in case of ABR.

2.7.4.2 Impact of Traffic Control Functions

Traffic control functions associated with ATM service categories can have a significant influence on the applicability of these categories. Two examples demonstrate these dependencies:

1. Decision between rt-VBR and CBR

Suppose a real-time application transmits with long bursts of high bit rate, but at a relatively low mean bit-rate. The burst tolerance of the source policing function is limited due to the burst sizes that can be handled by the network nodes. Consequently, it may be necessary to increase the sustainable cell rate far above the mean bit rate of the connection. Therefore, the customer's acceptability of the connection then depends on the pricing to determine whether the rt-VBR or a CBR service category is appropriate.

2. Application of UBR

For a multimedia application using UDP over AAL5, delivery of corrupted SDUs may be better than dropping whole AAL5 packets by a partial or early packet discard strategy. This is because partial information of the picture output may be extracted, and allow less corruption of the video signal seen at the application layer.

2.8 Summary

This chapter described the functionality of the Application Layer, AAL Layer and the ATM layer. The different methods of traffic control and differentiation between ATCs were stated. The difference between the ITU-T and ATM Forum standards was discussed and mapping strategies between ITU-T and ATM Forum QoS definitions formulated. Finally, a protocol stack selection strategy was

discussed that highlighted the importance of functionality, quality, price and traffic control.

The following chapter reviews only the ATM layer, as this is the most important layer for demarcation between users, service providers and operators. It also examines the methods to define, calculate and measure the network performance statistics that are needed to enable a network operator to provide performance guarantees.

3. ATM Performance Parameters

Performance parameters are used throughout the following chapters to obtain an understanding of the measurement events, places and methods. This chapter describes the parameters in detail. Network performance parameters are used to characterise the performance of an ATM layer VCC or VPC. Six network performance parameters have been defined, some indicate the user's network performance requirements, others are used by the network operator to measure the performance of the ATM network. One or more of these parameters may be offered on a per connection basis depending on whether the related performance objectives are supported by the network. The network can support different performance objectives by routing a connection through paths that have a similar performance objective. Alternatively, implementation-specific mechanisms within individual network elements can differentiate ATM connections leading to the supply of different Quality of Services.

The following network performance parameters can be negotiated:

- Maximum Cell Transfer Delay
- Peak-to-Peak Cell Delay Variation
- Cell Loss Ratio

The following network performance parameters are not negotiated:

- Cell Error Ratio
- Severely Errored Cell Block Ratio
- Cell Mis-insertion Rate

3.1 Network Performance Guarantees

An ATM network may support one or more performance objectives for each Quality of Service class. The negotiation of the network performance parameters is carried out between the end systems and the network(s) for each direction of the connection. The network will agree to meet or exceed the negotiated performance parameters for the duration of the connection as long as the end systems meet the

negotiated traffic contract. If the network cannot meet the requested parameters, the connection is rejected or, alternatively, a negotiation phase between the end system and the network is undertaken; the precise action being network specific.

The network performance commitments are probabilistic in nature and are only intended as an approximation of what the network expects to offer to the connection over its duration. The duration of a connection does not form a part of the traffic contract and thus the duration is unknown. The guarantees offered to the connection are dependent on the network statistics at the connection set-up. Therefore, the network performance may vary over the lifetime of the connection. Transient events, such as an interruption in the underlying transmission system, can cause a short-term deviation in the performance objectives of a particular connection. Performance guarantees can only be evaluated over the long term and over multiple connections with similar performance targets.

The mechanisms for transferring the network performance parameters through an ATM network can be found in [ATM010] and [ATM055]. These define how the network performance parameters may be negotiated between the end system and the network in quantified numeric units. These mechanisms supplement the ITU-T's "QoS Class" structure defined in [ATM056] and [I.356].

3.2 Cell Events and Measurement Points

The ITU-T has defined two events [I.356] that are important in establishing network performance. These events are:

- 1) Cell Exit Event. This occurs when the first bit of an ATM cell has completed transmission
 - out of an end system to a private ATM network across a "Private UNI" measurement point,
 - out of a private ATM network element into a public ATM network across a "Public UNI" measurement point,
 - out of an end system to a public ATM network element across a "Public UNI" measurement point.

2) Cell Entry Event. This occurs when the last bit of an ATM cell has completed transmission

- into an end system from a private network element across a “private UNI” measurement point,
- into a private ATM network element from a public ATM network element across a “Public UNI” measurement point,
- into an end system from a public ATM network element across a “Public UNI” measurement point.

The positions of these exit and entry points can be seen in figure 3.1. A combination of networks can exist between these measurement points; these have been highlighted within the figure.

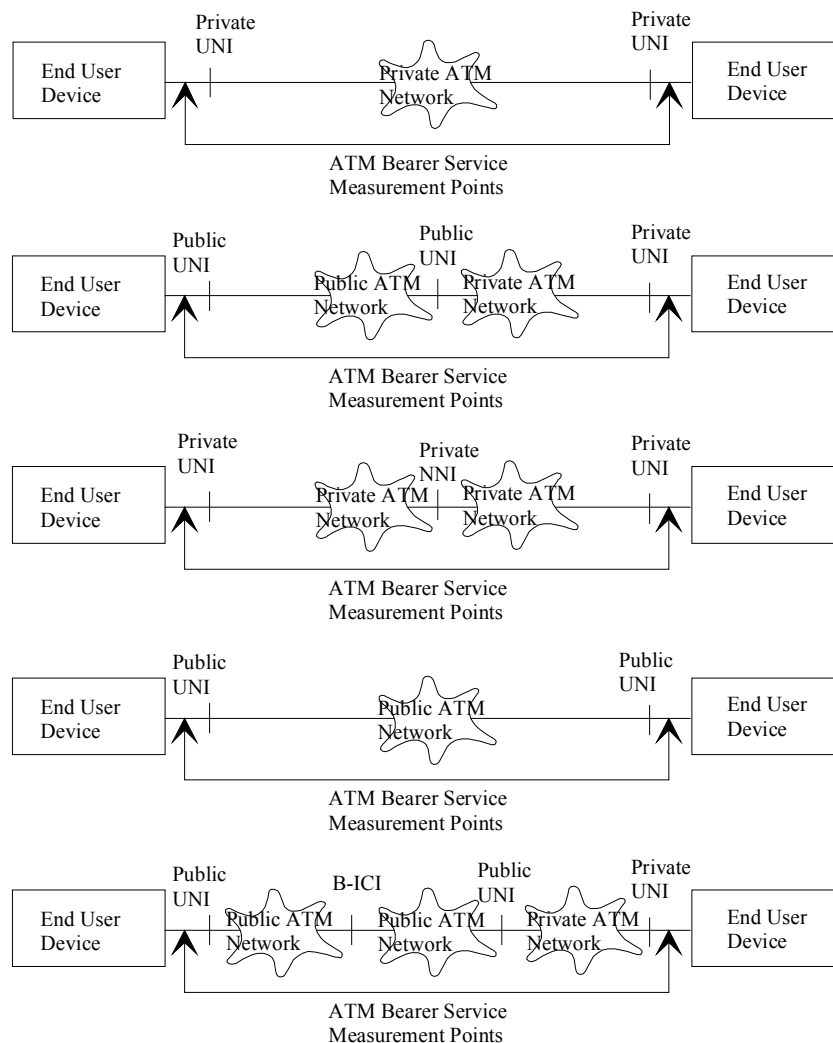


Figure 3.1: ATM Network performance reference points.

3.3 Cell Transfer Outcome

References [ATM056] and [I.356] define cell transfer outcomes between the above measurement points.

- **Successful Cell Transfer.** The cell is received within specified boundaries negotiated at connection set-up. The binary contents of the received ATM cell conforms exactly to the transmitted cell payload and the cell is received with a valid header field after validation by the header error check (HEC).
- **Errored Cell Outcome.** The cell has arrived within the negotiated bounds of the connection set-up. The binary content of the received cell differs from the corresponding transmitted cell payload or the HEC has determined an invalid header.
- **Lost Cell Outcome.** This is when no cell is received corresponding to a transmitted cell within the specified contracted parameters. This outcome applies to both “lost” cells and cells that arrive “late”.
- **Mis-inserted Cell Outcome.** This is when a cell is received for a connection for which there is no corresponding transmitted cell.
- **Severely Errored Cell Block Outcome.** When M or more “Lost Cell Outcomes”, “Mis-inserted Cell Outcomes” or “Errored Cell Outcomes” are observed in a received block of N cells transmitted consecutively on a given connection. M and N are defined in [I.610] and are explained in section 3.5.1.2.

3.4 Negotiated Network Performance Parameters

Using the measurement points in figure 3.1 network statistics can be obtained. These measurements are useful in gaining knowledge of the performance of the connection and, hence, provide an objective method for determining whether a connection has met the requested traffic contract. The measured network performance parameters will more than likely differ from the negotiated objectives for the connection at any given time. The “negotiated objective” is the worst case

of network performance that the network will allow; this includes periods when the network's utilisation is high. During periods when the network has a low utilisation factor, the negotiated objective may be significantly better than that specified in the traffic contract.

The measurement period of the network is very important. Transient events can cause measurements to appear much worse than the negotiated parameters if the measurement period is short. A measurement must at all times use a "significantly large number" to obtain valid results. This "significantly large number" is normally several orders of magnitude greater than the inverse of the desired ratio. For example, the $1 \cdot 10^{-4}$ CLR would require $1 \cdot 10^{+7}$ cells to be statistically significant.

3.4.1 Information Transfer Speed

The measured Cell Transfer Delay (CTD) is defined as the time a cell takes between a cell exit event at measurement point 1 (MP₁) and the corresponding cell entry event at measurement point 2 (MP₂). The CTD between the two measurement points is the sum of each link's transmission delay and each switch's processing delay along the virtual connection. Two end-to-end delay parameters are negotiated:

- Peak-to-peak CDV
- MaxCTD

In figure 3.2 is an illustrative probability density function of the CTD, which relates the peak-to-peak CDV to the maxCTD. The fixed delay includes the propagation through the physical media, delays caused by the transmission system, and fixed components of switch processing delay. CDV is caused by buffering and cell scheduling.

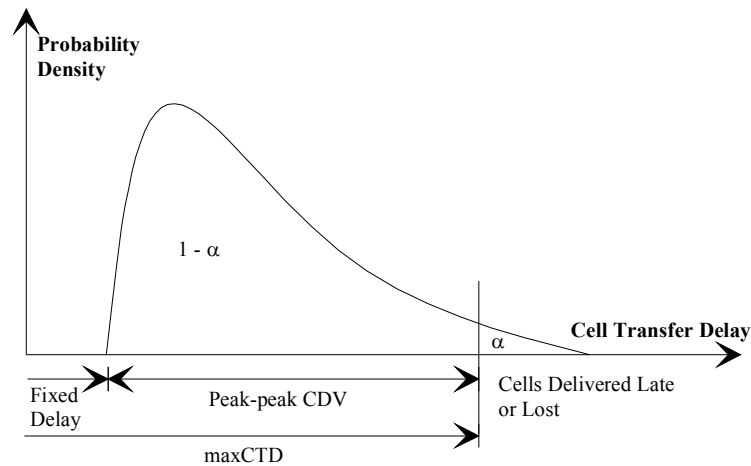


Figure 3.2: Illustrative Cell Transfer Delay Probability

An estimation of the upper bound of the speed parameters for variable length connections can be calculated using allocation rules as a function of “distance” and increasing “complexity”. This complexity consists of impairments that increase with additional switching stages and/or increases as more international or jurisdictional boundaries are crossed. That is, CTD is more a “distance” parameter, while CDV is a “complexity” one.

The connection is required to be divided into portions, as highlighted in I.356, and the “straight line” distance between the two boundaries of each portion should be known in order to calculate these expected worst case performance parameters. The three kinds of portions for transmission paths are defined in I.356. These are:

- National Portion (NP)
- International Transit Portion (ITP)
- International Inter-operator Portion (IIP)

Once the connection is divided into portions, the values of the upper bounds can be calculated. For every performance parameter, except the CDV, the end-to-end performance is the sum of the values in all of the portions. CDV accumulates as a function of the standard deviation of independent random variables.

Any connection that meets the performance objectives in section 8, or a connection portion that satisfies the allocation rules of section 9 of I.356 can be considered fully compliant with the recommendations.

There are standard methods to determine the mean cell delay and 2 pt CDV of a connection. The calculation is a dual process. First the straight-line distance must be found; from this value an estimate of the transmission path length is obtained. Finally, this is used to calculate the overall delay incorporating both the transmission path and switch delay.

3.4.1.1 Cell Transfer Delay (CTD)

Propagation delay, queueing, routing and switch delays affect the Cell Transfer Delay (CTD), these delays will vary in local and wide-area networks.

According to the ATM Forum the maximum Cell Transfer Delay (maxCTD) is specified as the $(1-\alpha)$ quantile of CTD, see figure 3.2. The CLR at the connection request time is used as an upper bound on α . As maxCTD and CDV is accumulated, a smaller value of α is selected at each switching stage so that maxCTD and CDV are over estimated.

The ATM Forum accumulates the CTD of each switching stage. In signalling, maxCTD is accumulated only in the forward direction (maxCTD_F). The CDV is accumulated in the forward and backward direction (CDV_F, CDV_B). The backward maxCTD is calculated as follows:

$$\text{MaxCTD}_B = \text{CDV}_B + \text{maxCTD}_F - \text{CDV}_F$$

This method of calculation is valid, as the fixed delay in the forward direction is the same as the fixed direction in the backward direction.

The ITU-T assumes that a network operator can define the performance criteria of the network connection by applying formulae to the straight-air-line distance between two end points. This calculation is less accurate, but gives a faster response to user requests. Using a formula-based approach, the user has a means to estimate valid network performance parameters for a particular connection.

Network operators prefer to withhold design details of the physical nature of their network. Thus, having a connection method that discloses accurately the delays between switches on a particular path opens the opportunity for a rival operator to determine the topology of a competitor's network. The ITU-T method to determine the CTD calculation is as follows:

For the *route length calculation*, D_{km} is the air-route (or straight-line distance) between two measurement points that bound a portion and R_{km} is the estimated Route Distance. Therefore:

If $D_{km} < 1000$ km, Then $R_{km} = 1.5 * D_{km}$

If $1000 \text{ km} \leq D_{km} \leq 1200$ km, Then $R_{km} = 1500$ km

If $D_{km} > 1200$ km, Then $R_{km} = 1.25 * D_{km}$

These calculations do not apply to portions with satellite hops.

To determine the *complexity calculation*, the ITU-T has defined a worst case value for the delay across a single ATM node to be 300 μ s. This is only for QoS class 1 services. The ITU-T has allocated a fixed number of nodes, N_{sw} , that can be crossed in each portion. Thus, knowing the number of portions in the connection, table 3.1 below can be used to assess the value of N_{sw} .

	National Portion	IIP(0)	IIP(1)	IIP(2)	IIP(3)	ITP
VCC	8 nodes	0 nodes	3 VP nodes	6 VP nodes	9 VP nodes	3 nodes
VPC	4 nodes	0 VP nodes	not applicable	not applicable	not applicable	3 VP nodes

Table 3.1: Number of switches assumed in each portion

When a connection portion contains a geostationary satellite hop, a delay of 320 ms is allocated for the CTD.

Finally, it is possible to obtain the *delay calculation* using the *route length calculation* and the *complexity calculation*.

$$\text{CTD (in } \mu\text{s)} \leq (\text{R}_{\text{km}} * 6.25) + (\text{N}_{\text{sw}} * 300)$$

The term $(\text{R}_{\text{km}} * 6.25)$ is an allowance for the “distance” within the portion, and $(\text{N}_{\text{sw}} * 300)$ is an allowance for the “complexity”.

3.4.1.2 Cell Delay Variation (CDV)

The CDV specification is important for CBR connection performance. Its value is necessary for dimensioning the elastic or play-out buffer required at the terminating end of the connection used to absorb the accumulated CDV. A maximum CDV value should be defined for a collection of network configurations. Crossing private and public networks with many switches in tandem should be incorporated in this calculation.

The ATM Forum defines the peak-to-peak CDV as the $(1-\alpha)$ quantile of the CTD minus the fixed CTD that can be experienced by each cell of the connection and during the lifetime of the connection. The term “peak-to-peak” refers to the difference between the worst case and the best case of the CTD. That is, the cell that equals the fixed CTD and the cell that has experienced the longest delay with the probability less than α . Assuming the fixed delay is the reference point for the two point CDV, then the distribution of the two-point CDV is the same as the peak-to-peak CDV. The control of the peak-to-peak CDV is difficult, therefore end-systems cannot negotiate arbitrarily small values of peak-to-peak CDV to meet jitter and wander tolerance [G.822], [G.823] and [G.824].

The ITU-T method to calculate the 2 point CDV for any portion supporting a QoS class 1 can be read from table 3.2. In this table, the general CDV measurements are assumed for national and international portions. The international portion can be further broken down into smaller classifications. The 2 point CDV accumulates as a function of the standard deviation of approximately independent random variables. Therefore, the total CDV of a single connection can be calculated by the following formulae:

$$CDV_{Total} = \sqrt{(CDV_A^2 + CDV_B^2 + \dots + CDV_N^2)}$$

Where CDV_{Total} is the CDV allocated for each portion of the connection.

	National	International				
CDV (ms)	1.5	1.5				
	National	IIP(0)	IIP(1)	IIP(2)	IIP(3)	ITP
CDV (ms)	1.5	0	0.7	0.9	1.1	0.7

Table 3.2: CDV Tolerance for Different Portions

The two point CDV describes the variability in the pattern of the cell arrival events. This is observed at the output of a connection at measurement point MP_2 , with reference to the corresponding cell stream at measurement point MP_1 .

The two point CDV (v_k) is defined as the variability of cell k , after passing through MP_1 and MP_2 , to the difference between the absolute transfer delay of cell k (x_k) and the reference cell transfer delay, this is defined as $d_{1,2}$. Therefore, the two point CDV of cell k is:

$$v_k = x_k - d_{1,2}$$

The reference cell transfer delay, $d_{1,2}$, is found through measurement and is equivalent to the smallest delay experienced by a cell in a cell stream on an unloaded network.

One point CDV describes the variability in the pattern of cell arrival events. One point CDV is observed at a single point with reference to the negotiated peak rate $1/\tau$, see [I.371].

The one point CDV for cell k (y_k) is measured by using a reference clock with an event arrival time of c_k and subtracting the actual arrival time (a_k) Thus $y_k = c_k - a_k$. The reference time c_k is calculated as follows:

$$c_0 = a_0$$

$$c_{k+1} = \begin{cases} a_k + T & \text{if } (c_k \leq a_k) \\ a_k + T & \text{otherwise} \end{cases}$$

By mapping the one-point CDV distribution, one can see both positive and negative values. The positive values represent cell clumping and the negative values represent gaps in the cell stream.

3.4.2 Information Transfer Dependability

One dependability parameter is negotiated:

- Cell Loss Ratio

This dependability parameter is applied to all service categories except the Unspecified Bit Rate ATC.

3.4.2.1 Cell Loss Ratio (CLR)

Cell Loss Ratio (CLR) is caused by most components of an ATM network. Queue capacities, reconfiguration of the transmission network, physical media errors all effect the CLR. Effectively, these conditions cause error to the header, buffer overflows or rejection of cells due to the non-ideal operation of the UPC mechanism.

Errors within the cell header at the physical layer are due to failures, protection switching and path reconfiguration. Buffering strategy and resources allocation can affect the CLR due to buffer overflow. Some networks provide large buffers and multiple levels of priority to obtain greater levels of efficiency from the transmission capacity. In this case, the CLR ratios tend to be better than a network not incorporating more complex buffering mechanism. The effect of CLR on higher layer protocols may also be distance dependent. For example, in a local area network higher layer protocols can request re-transmission of PDUs much quicker than if the connection was transmitted over a greater distance. Therefore, the requested CLR can vary for a similar service dependent on the locality of the end systems.

Multiple levels of delay priority and a combination of small and large buffers may be implemented to improve performance. When the CLP bit becomes significant,

the high priority cells will have a low CLR and the low priority cells will have a high CLR when the buffer becomes congested.

When nodes are connected in tandem, crossing traffic can influence the through traffic stream. When a greater amount of nodes are crossed, there exists a greater probability of cell loss.

Path reconfiguration from a longer to a shorter path also causes cell losses, because of the difference in propagation delay. Path reconfiguration in the transmission network is carried out in anticipation of different network requirements and out-of-service maintenance. An example of this is as follows: when a new path is set up, the traffic is transmitted on both routes until the new path is verified. Once the new path is operating correctly, the physical layer reconfiguration is made. This method is intended to reduce the interruption to users. Because of the difference in propagation delay, the longer path will hold a larger number of cells in transmission than the shorter path. Hence when the reconfiguration is made, some cells proportional to difference in path delay, will be lost.

The cell loss ratio is defined as:

$$CLR = \frac{\text{Lost Cells}}{\text{Total Transmitted Cells}}$$

Once cells have been included in a severely errored block, they will no longer form a part of the CLR calculation.

The negotiated CLR will last for the duration of a connection. However, the actual CLR may vary, as the value measured is dependent on the duration of the measurement and the number of cells sent within the connection.

The CLR objective will apply to the cell stream with the CLP=0 cell flow or the aggregated CLP=0+1 cell flow. This will be dependent on the particular ATC defined in the traffic contract.

3.5 Non-Negotiated Network Performance Parameters

The operator needs to determine the reliability of ATM equipment and the underlying transmission system, using these non-negotiated network performance parameters. These measures can be used by an operator to determine when maintenance is needed and where the faults are located.

3.5.1 Information Transfer Accuracy

There are three non-negotiated performance parameters. These performance parameters all deal with the accuracy of the information transferred. These parameters are:

- Cell Error Ratio
- Severely Errored Cell Block Ratio
- Cell Mis-insertion Rate

3.5.1.1 Cell Error Ratio (CER)

The cell error ratio (CER) is primarily caused by the error characteristics of the physical media. It is possible that this type of error is a function of distance and the characteristics of the media. Operational effects such as protection switching and network reconfiguration could also induce errors.

The CER is defined for a particular connection as: successfully transferred

$$CER = \frac{\text{Errored Cells}}{\text{Successfully Transferred Cells} + \text{Errored Cells}}$$

The successfully transferred and errored cells, within a severely errored cell block are excluded from the calculation for CER.

3.5.1.2 Severely Errored Cell Block Ratio (SECBR)

The severely errored cell block ratio (SECBR) is influenced by the error characteristics of the physical media and buffer overflows. The error

characteristics will probably increase as a function of distance and the characteristics of the media.

The SECBR is defined as following for a particular connection:

$$SECBR = \frac{\textit{Severely Errored Cell Blocks}}{\textit{Total Transmitted Cell Blocks}}$$

A cell block is a sequence of N cells transmitted consecutively. The size N is network dependent has been defined as either 128, 256, 512 or 1024, (additional block sizes may be defined in the future). A severely errored cell block is defined when M lost cells, errored cell or mis-inserted cell outcomes are present in a received cell block. The value M and in the implications of OAM flows on UPC and NPC mechanisms are for further study. A cell block is defined as a collection of user cells that lie consecutively between two OAM cells, see [I.610].

3.5.1.3 Cell Mis-insertion Rate (CMR)

The CMR is caused by undetected errors in the cell header. The HEC can correct a single bit or detect two errored bits. Therefore, if more bits are in error, then the cell could be interpreted as belonging to another connection. The errors will more than likely be caused by transmission error on either the underlying transmission network or physical medium. Once the header has been mis-interpreted a cell stream exists with the mis-interpreted cell header. These errors are likely to be rare and are dependent on the number of VPI/VCI values passing through a particular switch. As public switches have a greater likelihood of transferring a greater number of VPI/VCI cell streams, this type of error has a greater probability of occurring in public switches than private ones.

The CMR is defined as the following for a connection.

$$CMR = \frac{\textit{Misinsterted Cells}}{\textit{Time Interval}}$$

Cell mis-insertion is caused by undetected errors in the header and the cell being forwarded within the traffic stream of another connection. This parameter is a rate

rather than a ratio as the cause of mis-insertion is independent of the cell received in the corresponding connection.

3.6 Sources of Performance Degradation

Table 3.3 summarises the various methods of degradation on an ATM network, and how the degradation effects each network performance parameter.

	CTD	CDV	CLR	CER	SECBR	CMR
Propagation Delay	✓					
Bit Error Statistics			✓	✓	✓	✓
Switch Architecture	✓	✓	✓			
Buffer Capacity	✓	✓	✓		✓	
Traffic Load	✓	✓	✓			✓
Number of Nodes in Tandem	✓	✓	✓	✓	✓	✓
Resource Allocation	✓	✓	✓			
Failures			✓	✓	✓	

Table 3.3: Degradation of Network Performance.

The three most relevant to the work are the switch architecture, buffer capacity and resource allocation

- The switch architecture can significantly effect the network performance parameters. The switch matrix design, buffering strategy and the characteristics of the switch under load must be considered. The switch architecture can be blocking or non-blocking. The buffering capacity may use a single port or may be shared between multiple ports. The management of these buffers may incorporate a FIFO strategy or a more complex buffering strategy such as per VCC queueing. The switch matrix may also introduce errors under heavy load.
- The type and size of the buffering capacity has a great impact on the network performance parameters. For example, short buffers enable low latency

services, while large buffers allow statistical multiplexing of connections. This means that there can be significantly different performance characteristics.

- Depending on the CAC employed within the network, the resources allocated to the individual connections will vary. Therefore, as more or less network resources are allocated, the effect on network performance of each the connection passing through a particular link can change.

3.7 Summary

This chapter has reviewed only the ATM layer. The ATM layer is the most important layer as the UNI and NNI define the point of demarcation between users, service providers and operators. This chapter has examined the performance definitions, calculations and measurements that are needed to enable a network operator to provide performance guarantees.

The two previous chapters have outlined network performance and quality of service. There are many complexities in mapping between different standard bodies and selecting the most appropriate protocol stack. The following chapter describes a QoS mechanism that the author has proposed that can be provide QoS to even the most stringent applications in a simple, efficient, cost effective way.

4. Proposed QoS Framework

4.1 ATM Network Performance Parameters

As explained in the previous chapters, an ATM network has the important property that will allow different performance specifications over the same physical medium. The user of an ATM network is currently given a wide choice of ATM transfer capabilities (ATCs) to send different traffic types, i.e. CBR, rt-VBR, nrt-VBR, ABR and UBR. Network performance guarantees have been defined stringently for the CBR and rt-VBR, while nrt-VBR and ABR have only cell loss guarantees. Using combinations of ATCs (with and without priority mechanisms) gives the user the ability to change the QoS of the overall system. Network performance parameters have been used to define, dimension, and measure these ATCs. Control of these ATCs is possible through network control functions like UPC and CAC. Therefore, by varying one or more of these network performance parameters within an ATC, the QoS of the overall system can be changed.

The most stringent speed parameters of CTD and CDV have to be designed into ATM networks at the outset to allow CBR and rt-VBR traffic to be carried; stringent traffic categories cannot be applied as an afterthought. Once the most stringent services have been incorporated, then the addition of non-real-time services becomes a matter of adding the necessary control functions. CLR and CDV are inter-related. They are both influenced by network planning, control functions and the network's traffic mix. Therefore, the CDV and CLR for each service can be fine-tuned once an ATM network has been installed and knowledge of the traffic profiles and call patterns have been assessed. Thus, the CDV and CLR are "soft" parameters that can be adjusted according to users' performance demands and operators' utilisation requirements.

4.2 Quality of Service

Within Standards, the term "QoS Parameters" has been used liberally. QoS has been used to describe user-perceived subjective assessment on the application

layer, OSI interlayer performance achievements and ATM layer network performance. QoS within this thesis has been defined as the user's perceived subjective assessment. For example, a video broadcast with a low percentage of errors may present an acceptable QoS to the viewer; however, data transmitted to a computer with a small percentage of errors may not be acceptable. Therefore, the person who is most concerned about the overall system performance is also the person paying for the service. Hence, if the user perceived QoS is not met, the user will look to other networks for a better service.

The QoS is a subjective value the customer will use to label the service being offered. Subjective assessments such as “good picture quality”, “clicks every second or so” or “image freeze with rapid movements” can be used by the customer to describe the QoS. These parameters are not very useful to the network operator unless the reasons for the “clicks”, “freezing” and “losses” within the application can be identified in terms of network performance parameters, [WANG96]. From the work in this thesis, performance of end-systems and switches is measured to determine the appropriate levels of performance that should be supported by the network. Experiments with applications over these networks will allow operators to determine the effects of CTD, CDV and CLR on applications, allowing operators to identify the reasons for the “clicks”, “freezing” and “losses” that the customer may complain about.

In most cases QoS is not easily measurable: the effect of the performance impairment may not map directly to the effect on the user. For example, the effect of a lost cell on an MPEG video stream would depend on the Group of Pictures (GoP) structure, that is, the structure of Inter-coded (I), Bi-Directional (B) and Prediction (P) frames, [CHIA96]. If more B frames were present in the GoP then a single error would propagate for a longer time through the sequence of pictures, as complete picture refresh frames occur less frequently. Therefore, the end-to-end QoS in this case is very dependent on the higher layers and not just the ATM layer.

QoS cannot be guaranteed, but the resources necessary for QoS can. That is, the network operator can guarantee the “measurable” network performance

parameters: CTD, CDV and CLR. It is possible that end-to-end quality of service can be provided in collaboration with the equipment suppliers, but such an approach would be complex.

4.3 Customer, Operator and Service Provider.

There are many different QoS objectives to be met for the members of an ATM “community”. These members are: the customer; operator and service provider. Each member wants to be able to make gains (or profits) from the arrangement.

The customer wishes to transmit data quickly from A to B by the cheapest possible method with the highest QoS. Customers do not care if the transmission mechanism is Fast Ethernet, Frame Relay or ATM, although it would be beneficial if the operator could supply an interface that is compatible with the user’s existing interfaces. In fact, the network layer interface may not be the preferred interface to negotiate a contract, as the IP interface is becoming increasingly important with the accelerated growth of the Internet. This transport level protocol may take the main role in contract negotiation. Costs are an important factor not only to optimise towards the least cost, but there is also a need to have a cost system that is predictable, as customers do not want unexpected large bills. The customer wants a network that is able to improve the QoS or introduce the possibility of new advanced services. Finally and most importantly, the customer will look for a balance between pricing, performance and QoS.

The operator would like to benefit by maximising profits through selling network resources in an efficient manner. The operator could charge ATM users through a number of measurable criteria. Types of charging mechanism are not necessarily dependent on the service used; they may be solely based on the resources reserved or the resources used. This implies that an operator may be reduced to a “bit-carrier”. The charging scheme used may bill the customer according to a non-linear function. For example, low bit rate services may be charged at a high rate, medium capacity SVCs may be charged at a lower rate per bit, and high bandwidth PVC may return the lowest costs per bit.

The service provider could be the operator or an independent third party. The service provider will want to maximise the effectiveness of their service through standard marketing techniques. Again, as with all corporations, savings are desirable in communications costs, as the service provider will want to provide a cheaper service than rival companies, while maintaining profit margins. Like all parties involved, the service provider will want to minimise capital outlay and costs and increase revenues.

4.4 The New Approach

The thrust of the thesis has been to determine a method to provide end users with the required Quality of Service for their applications. There are many constraints, as outlined previously, involved in providing QoS and all must be considered before a solution is proposed. A particularly important aspect of provisioning QoS is a flexible and user-defined price/quality function. The QoS framework proposed will provide end-to-end QoS that is simple, efficient and cost effective.

The new approach this thesis proposes is to abolish the rt-VBR category, and in so doing remove the buffering requirements for real-time circuits from the network to the end-systems. The research supporting this thesis has been carried out through simulation, and experimentation on an ATM testbed, to prove the viability of the hypothesis.

4.4.1 The Problem

There are currently two real-time transfer capabilities; CBR and rt-VBR. These real-time categories are in direct “technological competition” with the existing circuit-switched networks, which transfer user information across the network as quickly as possible. To determine the performance of real ATM networks the underlying speed parameters of CTD and CDV are measured, and using these measurements, conclusions can be made about the degradation of the cell stream between end systems over an installed network.

Real-time CBR traffic is commonly used to transport services like video and telephony. These applications require good network performance and by

specifying the CBR category, the users will request the best performance from the ATM network. The CBR traffic category is easy for users to define and it is easy for the network operators to compute their load boundaries. The CBR applications in this category provide a constant stream of data that requires synchronisation between end-systems. This synchronisation is most important if good QoS for CBR applications is to be obtained. However, due to the nature of ATM networks, CDV is always added to the cell stream, which disturbs the end-to-end timing information.

Research reported in this thesis has formulated a scheme that is based on multiplexing CBR channels. When a network operator places traffic with a CBR characteristic over an ATM network it can coexist with traffic of the same type with minimal disturbance to network performance guarantees. This demonstrates that the CBR ATC, using a peak-rate allocation CAC policy, can be used to deliver the required end-to-end QoS.

The network does not only cause the degradation of speed parameters in CBR applications. Removing the CDV while reconstructing the constant rate data stream causes performance degradation. This additional “play-out” queueing required to remove CDV causes increased CTD, so causing the overall network performance to reduce and contributes to lowering the QoS as perceived by the users. By designing appropriate play-out queues and estimating the degradation of network performance, an assessment of CBR applications and the CBR category can be made.

CBR traffic is a relatively simple case; using a peak-rate allocation policy the traffic will cross the network with only a small amount of performance degradation. On the other hand, rt-VBR traffic is much more complex. To transmit rt-VBR traffic the user needs to know the characteristics of the traffic from the end-system in order to define the connection parameters. The user is expected to know the peak bit-rate, the sustainable bit-rate and the maximum burst size of the source. From this, the network operator uses the parameters to control and guarantee levels of performance across the network using policing and CAC. Figure 4.1 shows the PCR, SCR and the MBS being mapped into a conceptual

diagram of a dual leaky bucket UPC algorithm [I.356]. These control functions are necessary to protect the network from malicious users or failed equipment and are very abrupt in their functioning. If a user takes more than the allotted bandwidth, then the user's cells (depending on the ATC) will be discarded from the network, (or tagged for later discard). This is a good means of control from the users' perspective, only if the users know the characteristics of the traffic generated by the end-system.

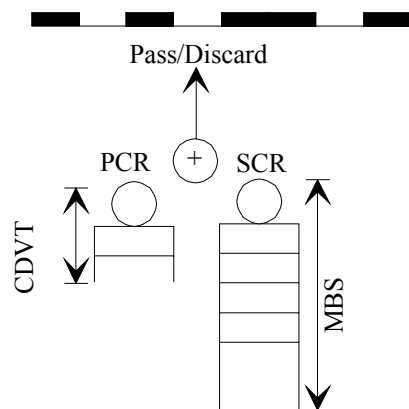


Figure 4.1: The Dual Bucket Policing Parameters

The users rarely know the values of the traffic descriptors that characterise their cell stream. During experiments conducted by the author on the ATM testbed, the characteristics of several *identical* applications were shown to be quite different. The main reason for the difference in traffic characteristics, found during a particular EXPERT “virtual classroom” event, was that one half of the classroom located in Switzerland was not as *vocal* as the other located in Canada. Hence, the quantity of traffic sent in opposite directions over the link, during a fixed period of time, were significantly different. Therefore, predicting the mean traffic rate is difficult and the specification of the traffic contract before transmission is, in most cases, impossible.

The operator has an equally demanding task computing the effects of different traffic types on the network. An operator uses the connection set-up parameters of PCR, SCR and MBS to determine the traffic characteristics, and the CTD, CDV and CLR for the required connection performance. Thus, to guarantee the performance of a connection, the operator selects a path according to the CTD and

CDV requirements, assesses the CLR across these paths, and increases the utilised bandwidth according to PCR, SCR and MBS.

To guarantee network performance for *every* rt-VBR connection on the network, the worst case traffic characteristics, (see figure 4.2) from a UPC mechanism must be modelled in order to ensure that the level of service offered to the customer is commensurate with the service delivered by the operator.

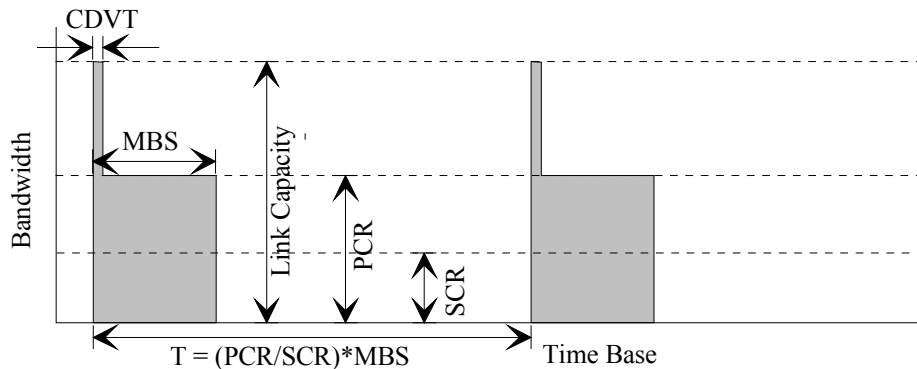


Figure 4.2: Worst Case Traffic from a Dual Leaky Bucket Mechanism

To demonstrate the difficulty in predicting the performance parameters of individual connections a simple simulation was set-up. Six different traffic types each requiring CLR better than $1 \cdot 10^{-4}$ have been defined in table 4.1.

	Inter-Arrival Time	ON Period in Cells	OFF Period in Cells
Type a	5	100	10000
Type b	5	40	1000
Type c	5	100	1000
Type d	5	10000	100000
Type e	5	10000	1000
Type f	5	CBR Connection	NA

Table 4.1: Six simulated worst-case traffic types.

To assess the overall performance of the network, a network operator may examine the CLR from each queue across the network as in figure 4.3. The six defined traffic types are multiplexed into a queue in figure 4.3, which resulted in a CLR of $2.086 \cdot 10^{-5}$. Under the operating definitions, the queuing element is working correctly.

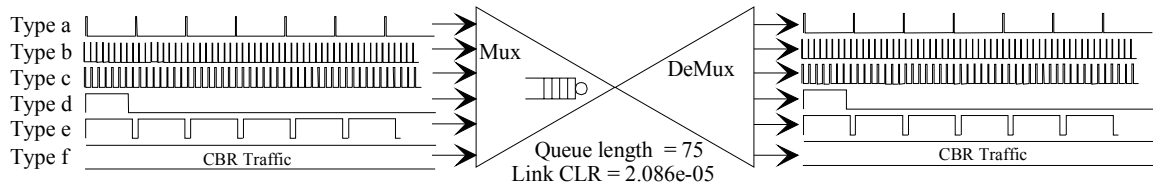


Figure 4.3: A switch multiplexing rt-VBR sources

However, when examining the individual connections passing through the multiplexer, connection of type “a”, in table 4.2, has a CLR of $2.034 \cdot 10^{-4}$ and because of its bursty nature is not obtaining the required network performance. In this case, due to the high CLR, the user will be charged for a connection that will not allow the end system to work. Under these circumstances the user may demand a refund or change to a different network operator.

	Cells Transmitted	Cell Lost in Multiplexer	CLR of each Connection
Type a	953613	194	$2.034 \cdot 10^{-4}$
Type b	3324030	194	$5.836 \cdot 10^{-5}$
Type c	6629450	247	$3.726 \cdot 10^{-5}$
Type d	6251168	221	$3.535 \cdot 10^{-5}$
Type e	19604825	183	$9.334 \cdot 10^{-6}$
Type f	20000000	145	$7.250 \cdot 10^{-6}$

Table 4.2: Traffic characteristics and network performance counters.

An operator could examine every connection at every node in order to determine that acceptable performance is being delivered. However, global knowledge of a network tends to be communicated slowly around the network and probably too slowly for corrective action to be taken.

4.4.2 The Solution

Previously, many analytical solutions have been proposed to the problem of multiplexing VBR traffic and “effective bandwidth” (EBW) is currently the most popular. There is a Poisson EBW Approximation [GRIF96],[GRIF97], Lindberger’s EBW [COST96], Large Buffer EBW [COST96] and Large Deviation EBW [KELL96]. These methods determine approximations of the resource

allocation necessary for each traffic type across the network, but cannot *guarantee* the network performance for every traffic type, particularly over real-time circuits. The use of large buffers and less stringent traffic contracts make these EBW methods more viable. The larger the average effective bandwidth, the lower the number of connections that are allowed on the network. A linear CAC function based on summing EBW can be used to determine the appropriate load boundaries.

These methods are complex and can be computationally heavy for the operator to implement, and also they assume that the users are able to describe their traffic characteristics appropriately before the connection begins. The simplest method to guarantee network performance would be to peak-rate allocate all rt-VBR connection across a network, treating all connections as CBR. In other words, abolish the rt-VBR category and only support CBR channels. The solution of peak-rate allocating connections, proposed within this thesis, may seem wasteful as no statistical multiplexing gain is sought within the network. It is, however, a method that will give guarantees of performance to *all* connections across the network and ensure that the service requested is rendered to the user. Traffic shaping allows the user to take some responsibility for network buffering through shifting switch queueing capacity to the end-system; this will enable the customer to make more efficient use of the requested CBR bandwidth and enable the user to select lower capacity CBR connections.

Customers tend not to consider other users of the network and will want to seize adequate bandwidth for their needs and more if possible. To prevent this, mechanisms are required that will ensure that customers share network resources in an optimum way. The mechanism proposed will provide a suitable service to the users by incorporating differentiation between ATCs, charging and shaping. The network operator uses a collection of ATCs that provide a required service to the user. For example, CBR is used for telephony and UBR for Internet type services. Within these service categories, the network operator can apply charges to place economic pressure on users to change their traffic characteristic to allow greater network utilisation. ATM network resources are a prime commodity,

particularly in the real-time categories. Economic pressure encourages customers to carefully assess their needs and will enable the operator to achieve greater utilisation by allowing more connections onto the network. Hence, the operator will increase revenue and this will lead to cost savings that could feed back to customers, who will get cheaper services.

Under economic pressure, the customer will want to change the traffic characteristics at the entry to the network. This can be done through shaping devices that add selective delay to the traffic stream. When traffic is clumped, the peak-rate of the connection is high; this clumping may only last for a short duration and spacing the cells evenly can make reductions in the instantaneous bandwidth. This “smoother” traffic enables the operator to increase the load boundary and hence the network utilisation. By shaping heavily, the VBR traffic stream can be smoothed to a near CBR profile. This will certainly allow the network to increase its multiplexing density and enable a greater number of users on the network. However, the shaper will increase the CDV of the application and could worsen the QoS or even stop the application from running. Therefore, a balance between shaping and QoS needs to be found so that the user obtains a suitable QoS and the operator increases the network utilisation. Shaping can be done to CBR connections to reduce the jitter within the network [LAN95] and to VBR connections to increase the nodal carrying capacity [ELWA97], [DANT94], and [NIES93]. The research in this thesis is orientated towards investigation of the trade-offs involved between QoS and traffic shaping.

QoS has been defined as the performance of the whole system. If an operator controlled the shaping devices, as proposed in [ELWA97], the operator would be unable to determine the price/QoS function required for a particular application. This is because the operator has no way to assess the QoS of the application as seen by the user. Therefore, it is more beneficial for the user to be responsible for the shaping device and hence the price/QoS function. In addition, the PCR of the proposed CBR-only connectivity is a simpler concept for the end user to master than the PCR, SCR and MBS combination of the rt-VBR ATC; it also allows simpler system configuration.

Shaping in end terminals means that the buffering capacity for the source is moved from the network to the user's end-system. This means that less buffering inside the network enables delay sensitive services to cross the network by the fastest possible means with minimal CDV degradation. This will enable competition between ATM and other circuit-switched networks. Customers who need a real-time service at a low cost can manipulate the price/QoS function by taking responsibility for buffering traffic in the end system.

4.4.3 The Research Methodology

The research methodology used within this thesis is experimentation and where this is not feasible, simulation. In chapter five, experiments are reported that researched the viability of the QoS framework hypothesis at a technical level. The research in this thesis covers the performance of ATM networks, the feasibility of high quality CBR services, QoS to network performance mapping, the user benefit of shaping while maintaining QoS, and increased network utilisation obtained by incorporating end-system shaping.

To be able to determine meaningful delay and CDV requirements, real ATM networks must be tested. CBR traffic and the proposed “shaped rt-VBR” stream will require a high-performance underlying network layer to achieve the demanding requests made from the service categories. Conventional CBR services require high performance circuits from the CBR ATC to compete with existing circuit-switched networks. The proposed “shaped rt-VBR” will be at the risk of being degraded by the shaper to a level where introducing additional CDV in the network could cause the application to fail. In experiment 1, real network testing is undertaken to determine whether the demands of the CBR service and the “shaped rt-VBR” service can be met, in order that the QoS framework proposed could exist over real ATM networks.

Thus, for synchronous systems a timing relationship between the end-systems is very important. Hence, using adaptive clock recovery and play-out buffering, CDV can be absorbed and timing information extracted. However, this additional play-out buffering process leads to worsening QoS objectives. In experiment 2,

this performance has been assessed to determine the additional delay and hence examine the degradation in performance of CBR services.

Having shown that an ATM network is able to deliver high performance circuits and the CBR service category can transfer cells with minimum performance degradation, experiment 3 examines the performance objectives from the application layer. In the QoS/Performance trade-off, bounded criteria must be met. An application needs a minimum level of performance to give suitable QoS, the minimum level depends on the type of application. A selection of real applications have been used with a panel deciding from a subjective point of view when an application ceases to give suitable QoS. When the application becomes irritating to the user, the values of CDV and CLR are recorded. This allows the minimum QoS of each application to be mapped onto network performance parameters.

Once the network and application bounds have been determined, research into shaping rt-VBR traffic can proceed. An initial study was undertaken to assess the types and characteristics of rt-VBR applications that are likely to be used across ATM networks. This study concluded that there are no native ATM applications available on the market and certainly there will not be any over the short and medium-term time scales. IP has a large user base and is still increasing rapidly. Services offered include standard data transferral e.g. FTP and WWW, as well as more interactive services, for example, Iphone, HQ Audio retrieval, Mbone Picture and Sound distribution and Multimedia terminals using low picture rates.

IP is the most important protocol for data communications, and this type of traffic will be significant throughout this thesis, as VBR traffic will use an IP over ATM model. The fact that there are no native ATM applications on the market means that bursty traffic becomes a predominant feature of the traffic characteristic on the ATM networks due to IP layers. In experiment 4, traffic from a real application is used with shaping devices to determine the delay performance in the end-system. Hence, the experiment examines the feasibility of shapers for reducing the PCR parameter, without losing the applications QoS. This allows the customers to obtain a cheaper service for the same application QoS. Finally,

experiment 5 examines the consequence of shaping in terms of network utilisation. The multiplexing density of a collection of homogeneous multimedia terminals is assessed by obtaining the admissible load boundary. From this, the network operator can share the benefits of having traffic on the network that enhances a higher multiplexing density by incorporating shaping in the end-system.

5. An Experimental Approach to ATM QoS

Five experiments have been defined to determine the validity of the hypothesis. Experiment 1 examines the performance of real-time ATM networks on different network scales. Experiment 2 examines CBR services over real-time networks and examines principally the effects of CDV on the destination terminals. From this, it is possible to determine the ability of ATM network to deliver stringent network performance. Experiment 3 determines the degradation of the application layer QoS with different performance conditions. Experiment 4 investigates the robustness of rt-VBR application to peak rate traffic shaping. This allows rt-VBR traffic to efficiently use lower capacity CBR channels to cross the network with the required QoS. Finally, experiment 5 researches the benefit of traffic shaping rt-VBR sources on network utilisation.

5.1 Experiment 1, Cell Transfer Delay and Cell Delay Variation Measurements

5.1.1 Introduction

ATM networks can accept parameters that specify the requirements for particular connection characteristics. The ATM transfer capability and the respective negotiated performance parameters are shown in table 5.1.

	CBR	rt-VBR	nrt-VBR	ABR	UBR	ABT
CTD	Specified	Specified	Unspecified	Unspecified	Unspecified	N/A
CDV	Specified	Specified	Unspecified	Unspecified	Unspecified	N/A
CLR	Specified	Specified	Specified	Specified	Unspecified	N/A

Table 5.1: The ATCs with their respective performance parameters

To achieve end-to-end QoS a user will have to specify meaningful delay and delay variation parameters in the connection set-up. The objective of Experiment 1 was to determine the performance of an ATM network by measuring switching devices. This provided a comparison between the performance requirements of applications and the performance capability of the network. There has been little research on the performance of real-time ATM networks, although [BANE97], [DASI97] and [SRIK96] do discuss the performance of the vBNS ATM network

in North America, but do not compare the delay to standard methods of calculation. The author contributed to a publication on network performance in [BRIN96]. Experiment 1 develops from [BRIN96], which examines performance on a switch-by-switch basis, to provided measurements bounding the maximum network performance parameters on multi-switch networks. Once the underlying network performance was bounded, further modification to services, detailed in later experiments, was justified.

5.1.2 Experiment Set-up

5.1.2.1 Individual Switching Elements

ITU-T recommendation I.356 [I.356] states that each switch should have no delay greater than 300 μ s for real-time circuits. A number of ATM switches were examined to determine the maximum and minimum delay. The minimum delay occurs at low buffer occupancy where the predominant delay is due to cell processing; the maximum delay is obtained when a cell crosses the switch that has a high buffer occupancy. Thus the total delay is the sum of the cell processing delay and the maximum queueing delay. The test configuration shown in figure 5.1 is used to determine ATM switch delay performance. The test traffic had low bandwidth and contributed little to switch loading effects, allowing the highest switch performance to be measured. Background traffic was introduced later; this caused the buffer occupancy to rise to the maximum. Hence, measurements of a switch under heavy load were obtained.

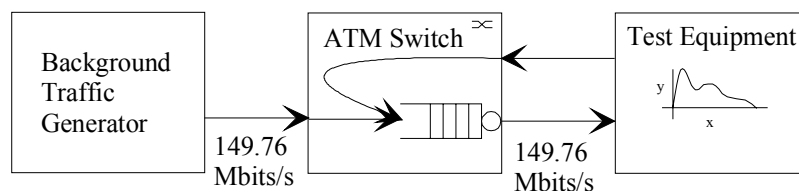


Figure 5.1: Hardware configuration to test a single network element.

The transmission delays between switches were generally negligible. The delay on an optical fibre is fixed when considering ATM cell transfer; photonically it may vary slightly. The normal value taken to judge the speed through a fibre is

approximately 5 μ s/km. The results from the experiments are considered in section 5.1.3.

5.1.2.2 Network of Switching Elements.

Performance measurements were made over a network of switches, which allowed the observation of network characteristics in two connections of different lengths. CBR test traffic was generated by the test equipment and sent through these connections. From the 1.6 million cells of each connection, it is possible to evaluate the CTD and the CDV. The values obtained were then compared with ITU-T recommended upper bounds in order to verify whether the real network complied with these recommendations. The calculation of the speed parameters of QoS 1 (stringent class) from the ITU-T recommendation was selected in order to compare against the experimental results. This QoS class is selected, as this is the only ATC class that imposes CTD and 2-pt CDV performance objectives.

The performance objectives for CTD are upper bounds on the underlying mean CTD of the connection. Although, individual cells may have transfer delays that exceed this bound, the average CTD for the lifetime of the connection (a statistical estimator of the mean) should normally be less than the CTD bound. The performance objectives for 2-point CDV are upper bounds on the difference between the 10^{-8} and the $1-10^{-8}$ quantiles of the underlying CTD cumulative distribution for the connection.

The two connections measured were set-up over international and intercontinental ATM networks. One connection was set-up across Europe, and the other involved a transatlantic submarine cable to Canada. The European connection was a peak rate allocated VC connection that originated from the testbed in Basel (Switzerland) and went via Zurich, Cologne (Germany), finally terminating in Leidschendam (Holland), see figure 5.2. The intercontinental connection was a peak-rate allocated VC connection starting at the testbed in Basel (Switzerland), through the European ATM network, then on to Canada via a submarine cable, see figure 5.3. Both the connections can be divided (following I.356) into separate portions, and an estimated value of the CTD and CDV obtained.

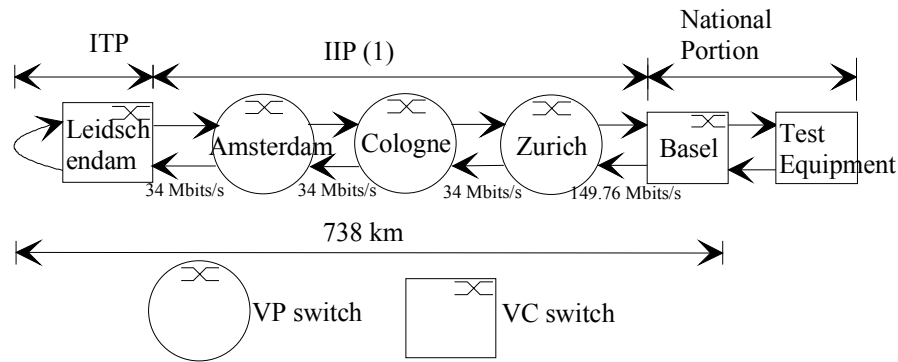


Figure 5.2: Basel-Leidschendam-Basel Portioned

The international connection had two National Portions, two International Inter-operator Portions (IIP(1)) and one International Transit Portion (ITP).

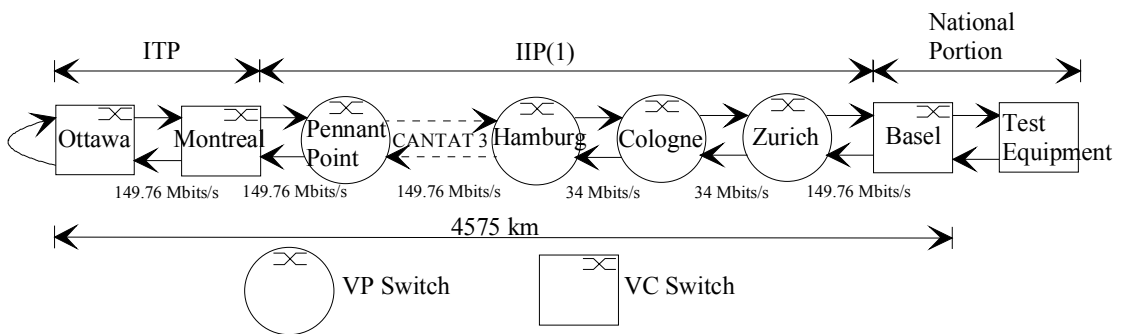


Figure 5.3: Basel-Ottawa-Basel Portioned

The intercontinental connection had two National Portions, two International Inter-operator Portions (IIP(1)) and one International Transit Portion (ITP).

5.1.3 Results

5.1.3.1 Single Switch Element

The minimum and maximum cell transfer delays were used to determine the best and worst delay performance of each switch. From these delay measurements, an estimate was calculated of the queue length within the switch. This measurement, see table 5.2, was necessary as manufacturers were unwilling to release information on buffer sizes.

Switch	Min Delay	Max Delay	Estimated Queue Length
AT&T, RUM	57.25 μs	140 μs	30 cells
Philips, LaTeX (Pure ATM)	47.03 μs	188.80 μs	52 cells
Philips, LaTeX (SDH)	151.31 μs	314.89 μs	61 cells
Alcatel, ALEX	127.46 μs	498.24 μs	140 cells
Ascom, AAU	104.28 μs	769.51 μs	254 cells
Fore, ASX-200	68.84 μs	766.78 μs	293 cells

Table 5.2: Switch Delay Measurements

Some of the switches have a maximum delay greater than 300 μs . However, the three switches well over this value have user defined maximum queue lengths and so the maximum delay in the switch can be reduced below 300 μs . Those switches with larger buffers allow the output buffer to have different scheduling mechanisms and hence incorporate both “real time” and “non-real time” circuits. This gives the switch a greater flexibility to handle different classes of traffic, see for example [ASX295], [SBA94] and [BRIE98].

Unlike the ITU-T, the ATM Forum expects the network operator to test each switch and link when the network is installed, so that connection set-up procedures can use the link-by-link performance parameters. When a connection request is made, each set of performance parameters along the connection was collated to determine if a connection could be accepted or not.

For example, in figure 5.4, a set of switches was defined along an end-to-end path of an ATM connection. Applying the results obtained in table 5.2 the end-to-end delay path was established according to the ATM Forum recommended calculation.

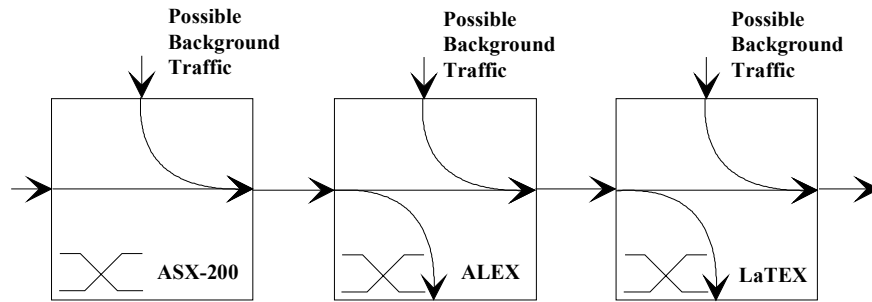


Figure 5.4: Example network.

Making the assumption that the cable delay was negligible between each switch, then the data in (table 5.3) can be considered a maximum, or worst case estimate, for the CTD and CDV in the forward and backward direction, according to the ATM Forum methodology.

The α quantile needed for the CTD and CDV calculation has been defined as the queueing capacity of the switch. This means that the real value of α has been overestimated to provide worst case results.

Switch Label	MaxCTD	pk-to-pk CDV	MaxCTD _F	pk-to-pk CDV _F	MaxCTD _B	pk-to-pk CDV _B
ASX-200	769.51 μ s	700.67 μ s	769.51 μ s	700.67 μ s	1456.55 μ s	1213.22 μ s
ALEX	498.24 μ s	370.78 μ s	1267.75 μ s	1071.45 μ s	687.04 μ s	512.55 μ s
LaTEX	188.80 μ s	141.77 μ s	1456.55 μ s	1213.22 μ s	188.80 μ s	141.77 μ s
		Total:	1456.55 μ s	1213.22 μ s	1456.55 μ s	1213.22 μ s

Table 5.3: Estimation of Delay in an Example ATM Network.

The ATM forum method is an approximate way to find the CTD and CDV through a network of switches. It does not determine the CLR, the mean delay, or the IAT p.d.f of the arriving cell stream. This is because these parameters are dependent on traffic characteristics of all the sources using the network. To be able to determine these other network performance measures, a greater knowledge of the traffic on the network is needed. A more detailed method for obtaining the delays through a network of switches is offered in [DEL48] and [VLEE95], where the delay distribution of a series of switches is convolved to produce the end-to-end delay.

From the results in table 5.3, the values of CDV seem to be very large. The fixed delay of this network is 243.33 μ s and the peak-to-peak CDV contribution is 1213.22 μ s. The actual measured CDV would be based on the resource allocation policy and the type of traffic on the network. The CDV value can be substantially reduced if a CAC policy of peak rate allocation were enforced. Cell scale queueing is predominant in peak rate allocated networks and Poissonian analysis is one method to approximate the maximum queue lengths. Hence, for a network load of 80%, a switch queue length of 40 cells would be sufficient for a cell loss probability of $1 \cdot 10^{-8}$. Thus, using the method of Poisson analysis, the maxCTD would be 570.49 μ s and a peak-to-peak CDV value of 327.16 μ s for the network of switches. These values are better approximations to the expected CTD and CDV of real-time services as highlighted in table 5.3.

5.1.3.2 Network of Switching Elements

From the ITU-T standards, the expected upper bounds for the CTD and the 2-point CDV of the *international connection*, (as defined in I.356 for each portion) are shown in table 5.4.

	EUROPEAN	
PORTION	CTD (ms)	CDV (ms)
National	2.4	1.5
IIP(1)	7.82	0.7
ITP	0.9	0.7
TOTAL	21.34	2.44

Table 5.4: Expected Delay and CDV measurements

The probability density function of the CTD over the connection from Basel to Leidschendam is shown in figure 5.5. The measured mean transfer delay calculated from the p.d.f. is 17.636 ms with a standard deviation of 0.015. The measured value is lower than the value recommended by the ITU-T of 21.34 ms. From this result, it can be seen that the measured connection is better than the delay approximation defined by recommendation I.356.

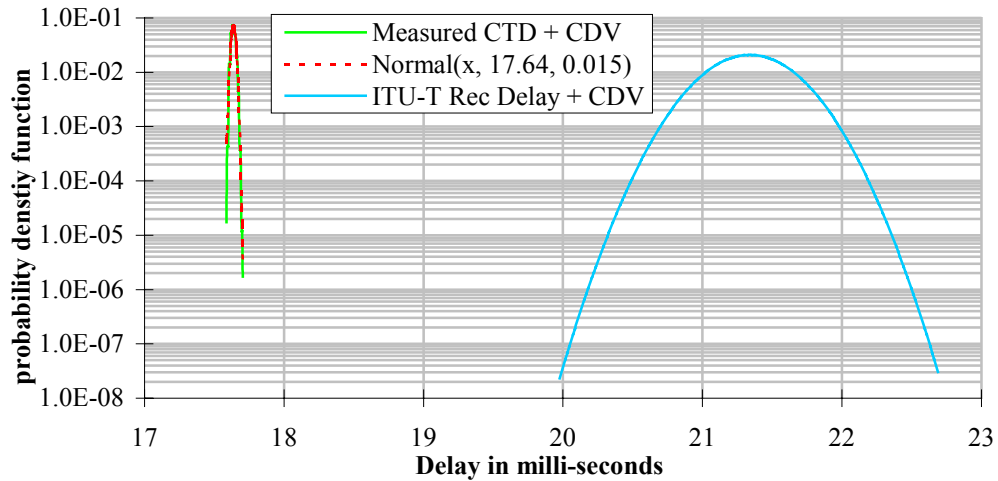


Figure 5.5: Cell Transfer Delay Basel-Leidschendam-Basel

The upper bound of the 2-point CDV recommended by the ITU-T applies to the difference between the $1 \cdot 10^{-8}$ and the $1 - 1 \cdot 10^{-8}$ quantiles of the underlying CTD distribution for the connection. The difference between the measured minimum and maximum CTD is used instead, as only 1.6 million cells are captured. Therefore, the maximum 2 point CDV is 0.12 ms. This value is far from the upper bound recommended by the ITU of 2.44 ms. The Pan-European network is a peak-rate allocated network and the only type of queueing is cell scale queueing. The large difference between the measured CDV and the calculated CDV is caused by the queue distribution of each switch not deviating far from the empty state. This results from the small number of independent sources, a peak-rate allocation policy and a low network utilisation. Thus, the measured CDV result is much smaller than expected.

The expected upper bounds for the CTD and the 2-point CDV for the *intercontinental connection*, as defined in I.356 for each portion are in table 5.5.

	TRANSATLANTIC	
PORTION	CTD (ms)	CDV (ms)
National	2.4	1.5
IIP(1)	35.33	0.7
ITP	3.99	0.7
TOTAL	79.45	2.44

Table 5.5: Expected Delay and CDV measurements

The measured probability density function for the delay of the connection from Basel to Ottawa can be seen in figure 5.6. The mean CTD for the round-trip delay is 109.54 ms and has a standard deviation of 0.0153. The mean CTD is larger than the value recommended by the ITU of 79.45 ms. This means that the connection is not compliant with the CTD recommended by the ITU-T. Within the IIP(1), a submarine cable with long delay is incorporated. This seems to have caused an underestimation of the calculated maximum delay, therefore enhancements to the calculation method are needed for long repeated spans.

The difference between the minimum and maximum CTD is 0.17 ms. This CDV is again much lower than the recommended value of 2.44 ms. The reason for the low CDV, as highlighted earlier, is the small number of independent sources, peak-rate allocation policy, and a low network utilisation

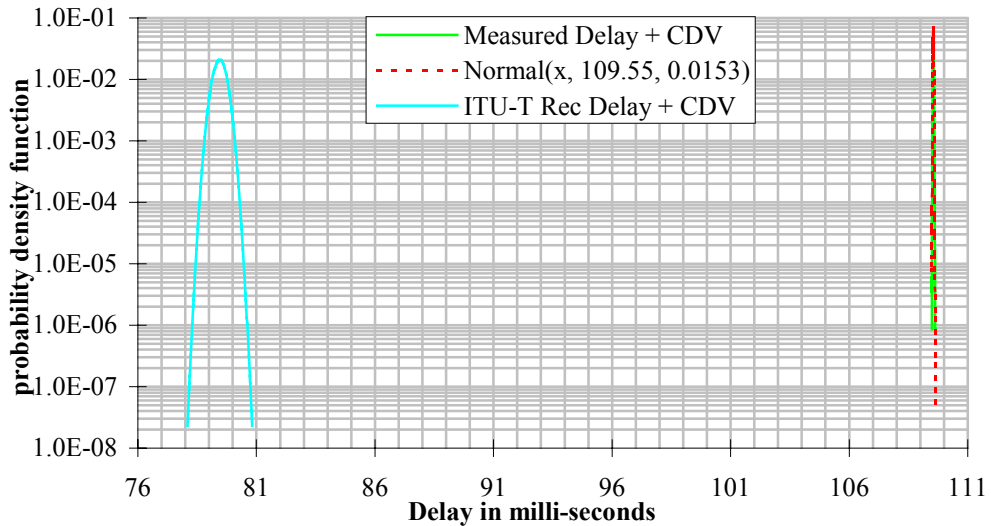


Figure 5.6: Cell Transfer Delay Basel-Ottawa-Basel

5.1.4 Summary

Experiment 1 has provided a greater understanding of the expected speed performance of ATM switches. Table 5.2 gives some indication of how much a network's performance can change under load, and estimates of the range of delay. These results present a method to determine the expected performance of small LANs by using the ATM Forum's recommended calculation for speed parameters.

When measuring the delay through a cascade of ATM switches over distance, the ITU-T recommendation was a rough approximation for switches designed to carry the stringent class of traffic. The measured delay across the Pan-European connection was less than the recommended values and had enhanced performance. However, there was a discrepancy in the results of the intercontinental link and revisions in the algorithm to determine speed parameters over long cable lengths, (such as transatlantic submarine cables) are needed.

The two point CDV values in both the long distance connections were much lower than the recommended values, and both the measured values were quite similar. This was due to the low number of independent sources and the low network utilisation. On examination of the measured traffic it was possible to determine

that the CDV did not depend on the length of the connection, only on the “complexity”.

Experiment 1 has bounded the underlying network speed performance of ATM networks. These measurements will be used as a benchmark in the following experiments to determine whether an ATM network can meet the performance required for stringent services from end-to-end, while maintaining user Quality of Service.

5.2 Experiment 2, Investigation of End-System Speed Degradation by CBR Play-Out Queueing

5.2.1 Introduction

One of the functions of ATM networks is to transport constant bit rate services using AAL1; this is referred to as Circuit Emulation, [ATM032]. During the B-ISDN introduction phase, existing telephony applications will require inter-working units placed in the DTE, CPN or at the LEX to convert from the “older” communication protocols to a form that can be transmitted over B-ISDN networks. These constant bit rate sources present a particular challenge to ATM networks, as an accurate timing relation has to be maintained between the source and destination.

When a source generates an AAL1 traffic stream over an ATM network, the sink cannot expect to receive exactly the same traffic characteristics as was transmitted. This is due to the statistical time division multiplexing used in ATM networks, [LOUI94]. AAL1 type services will be transported over higher priority “real time circuits”. AAL1 will use time priority and hence have a separate small buffer that will be served before other lower priority buffers. Employing a priority mechanism allows real time traffic to cross the network in a near deterministic fashion and increases the utilisation of the network, by identifying lower priority traffic that is not CDV sensitive, [NASE96], [BOLL96], [LAND94] and [LAND95].

To provide a good QoS for CBR applications, good network layer performance is necessary. CBR services are time critical and CTD needs to be minimised. However, as cells are statistically multiplexed, the inter-cell time can be disturbed, causing CDV. To compensate for CDV, the CBR end-terminal needs to add more delay. This can have a significant effect on the QoS of the application, as the end-to-end delay needs to be minimised to reduce annoying silence periods in an interactive conversation. However, the buffering must be sufficient to absorb the CDV effects that can cause loss of sound. The aim of this experiment is to determine the optimal queue size in the end terminal. From this, the degradation of the end-to-end performance was assessed for a CBR connection accounting for the additional end system functionality.

5.2.2 Play-Out Buffering

The QoS of an application is very dependent on the functionality and dimension of the play-out buffer located in the receiving terminal equipment. This buffer is investigated within Experiment 2 to determine the degradation of speed due to the inclusion of this buffer. The play-out buffer primarily consists of a queue and server that is used to remove the cell delay variation at the expense of additional delay.

5.2.2.1 Queue Functionality

In figure 5.7, the cells cross the network and the CDV variation is superimposed on the cell stream by passing the traffic through a cascade of queues. To reconstruct the original CBR traffic profile, additional delay is added to selected cells, nominally to bound the worst delayed cells. Thus, one can determine a fixed delay boundary that will remove all, or most, of the CDV. There exists a problem with the buffering mechanism, as a half-filled queue is needed to absorb the variation in delay. If the queue is too short then the CDV caused by the network will not be smoothed or losses in the traffic stream would result. If however, the queue is too long then the additional delay will reduce the maxCTD performance of the network.

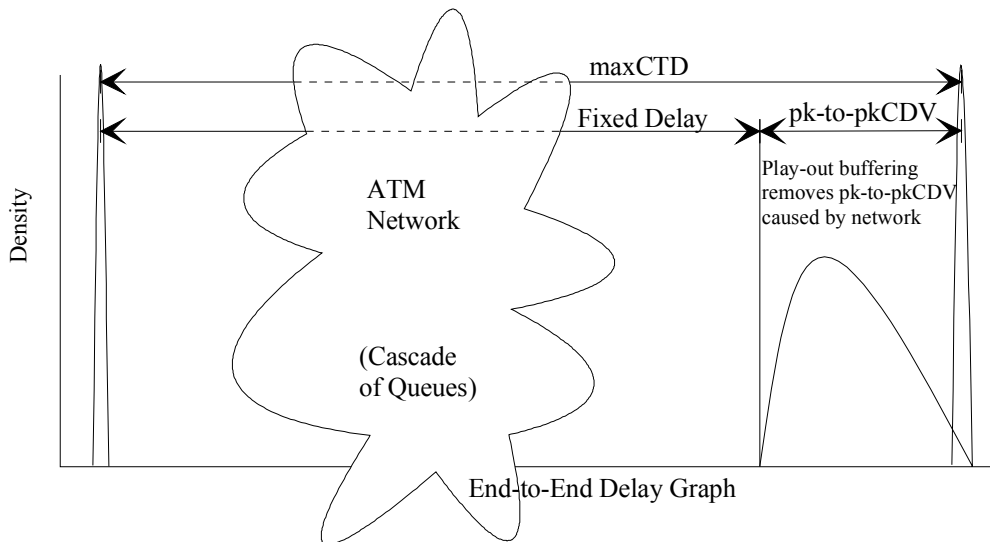


Figure 5.7: Reconstructing the traffic profile with a play-out buffer

5.2.2.2 Server Functionality

The server within the buffer is designed such that it serves the queue deterministically at the mean arrival rate. However, one problem with ATM circuit emulation is that no clock information is transferred and adaptive clock regenerative synchronisation is used in the end equipment. Thus, it is highly probable that the receiving clock will not be synchronised to exactly the same frequency as the transmitting clock. The receiving clock will initially make a “best estimation” of the mean arrival rate. The queue level of the play-out buffer is then used in an adaptive manner to control an oscillator that will advance or retard the receiver’s application clock (and server rate). From [DEL31] (quoted below) it has been shown that this technique is employed to recover a “best estimate” of the transmitting system’s clock of the high quality audio terminal. The following passage has been taken from [DEL31].

“Recovery of the continuous bit stream

The recovery of the continuous bit stream is a special problem with regard to transmission of digital signals in the ATM network. At the receiving end, no reference clock is available, since the data are transmitted in the form of cells and thus do not arrive continuously. Another problem is that the gaps between the cells of a virtual connection are not of a constant length. The reason for this is

that the data at the transmitting end are stored in a buffer until an information field is filled in. The cell, however, can only be released in certain time intervals - to replace idle cells in the cell stream - which leads to variable waiting times. On the other hand, the cells are subject to variable delays in the exchanges, since free time positions in the cell stream are used in the outgoing direction and therefore the cells have to remain in the buffer for different periods of time. All this leads to a considerable variation of cell delay named variable cell delay.

It is now necessary to recover a continuous data stream corresponding to the original one from this highly discontinuous data stream by avoiding information losses. This is quite easy by applying an adaptive method for clock recovery. In doing so, a buffer required for avoiding information losses is simultaneously used for the control of a clock oscillator. If a low-pass characteristic of a sufficiently low limiting frequency (mHz range) is assigned to this loop, a digital signal with a very low frequency clock jitter is obtained and can be easily processed by the terminals.”

When CDV is added to the cell stream, the queue occupancy in the terminating equipment will also deviate by a proportional amount. As the CDV increases, the deviation increases in the queue of the play-out buffer until it reaches either, the “empty” state or the “overflow” state. In the empty state, the play-out buffer will “reset” itself and wait until enough cells have arrived to half fill the queue; conversely, when the queue is in the “overflow” state, the buffer will “flush” and again the buffer will wait until the mid-position is entered in the queue before playing out any cells.

5.2.3 Experimental Set-up

Experiment 2 used a simulator constructed by the author in the SIMULA programming language, as to access and vary buffer dimensions in real ATM end-systems would be impractical. Simulations were used to assess the queue length requirements of a number of bandwidths requiring an AAL1 circuit emulation service. This was to determine the optimum size of a play-out buffer with respect to the bit rate of an application.

When a CBR data stream is transferred over an ATM network, the input stream is mainly disturbed by the background load on the network. This disturbance is an unknown quantity as large ATM networks with many users are not available for measurement. A worst case model for a CBR source queuing in a large ATM network is a Poisson process. This is a reasonable estimate for “real time” circuits as the queues within the network are designed to absorb only cell scale queuing and bursty traffic will either be peak-rate allocated (or something very close to peak-rate allocation).

Sets of simulations were performed that used a Poisson arrival process as a CBR source that fed into a model of a play-out buffer, see figure 5.8 for the set-up configuration. The service rate of the buffer was set to the mean rate of the input traffic. As the mean of the traffic was already accurately known, the server had no advance or retard function, (as is needed in real hardware).

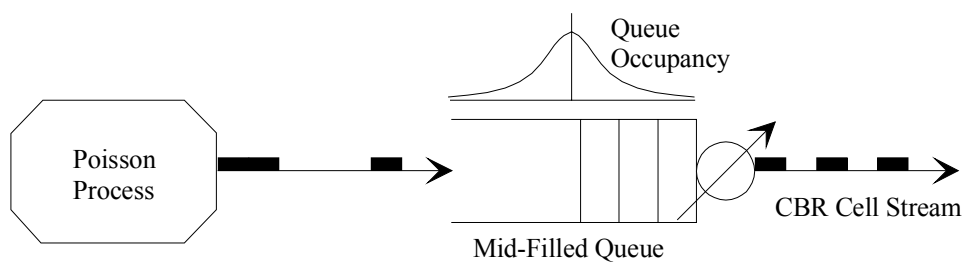


Figure 5.8: Play-Out Queue Experiment Set-up

5.2.4 Results

The simulated queue distribution of the play-out buffer was assessed. Figure 5.9 shows two results using the Poisson model, one with a mean bandwidth of 210 kbit/s (IAT = 740) and one with a mean bandwidth of 38.88 Mbit/s (IAT = 4). Both a Poisson and Normal Distribution were used to describe the queue occupancy and it became apparent that the play-out buffer distribution could be described best using a Normal Distribution.

Figure 5.9, shows two examples of play-out queue occupancy. From the graph, the requirements in terms of queue length for low bit rate connections were

smaller than for the higher bit rate traffic. This feature was repeated throughout the set of simulations to form experiment 2.

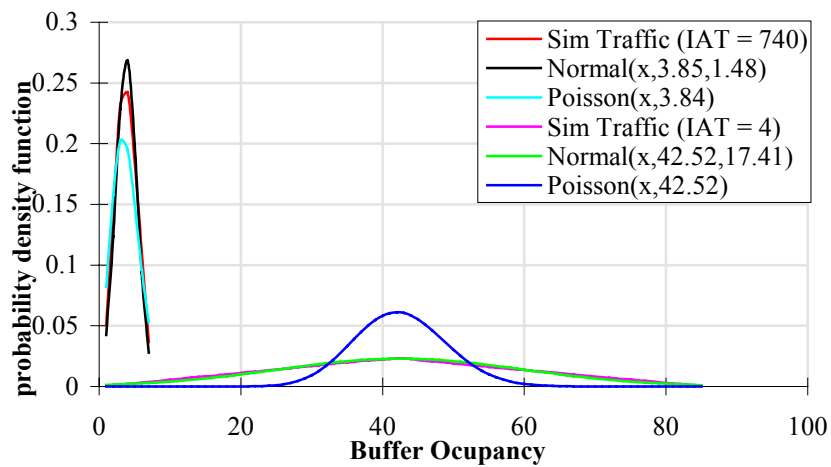


Figure 5.9: Queue Occupancy for 0.210 & 38.88 Mbit/s CBR sources

Using the assumption that as CBR traffic crosses an ATM network, it tends to a Poisson arrival process, it is possible to relate the known arrival rate to an optimum queue length for the play-out buffer. This prediction can be made by mapping the standard deviation of the arrival rate to the standard deviation of the queue distribution.

Thus, in figure 5.10, σ_p is the standard deviation of the input traffic. As the input traffic is a Poisson process then: $\sigma_p = \sqrt{\text{Mean Interarrival Time}}$. σ_n is the standard deviation for the queue occupancy and thus one can relate the input arrival rate to the required buffer size.

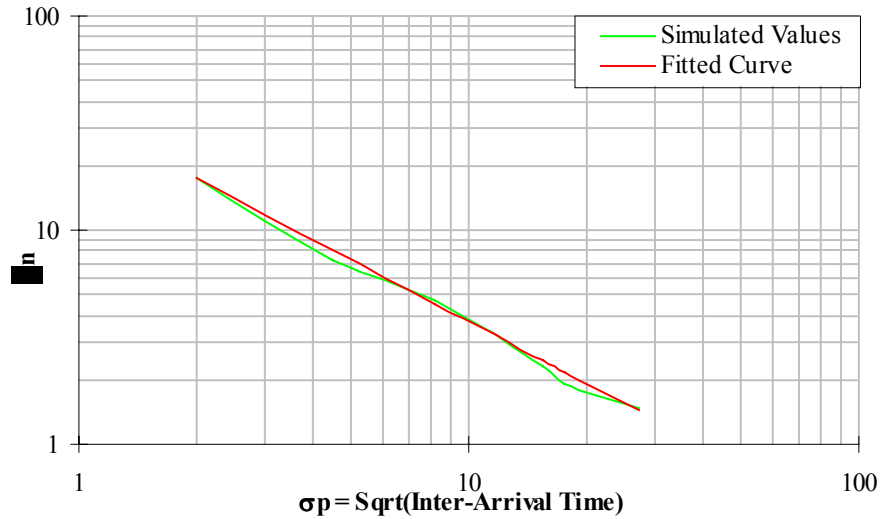


Figure 5.10: Mapping σ_p of the Arrival Process to σ_n of the Queue p.d.f.

Analysing the results, a non-linear function can be formulated that gives the following equation:

$$\sigma_n = f(\sigma_p) = \sigma_p^m \cdot 10^c \quad \text{if } 2.0 \leq \sigma_p \leq 27.2$$

Where:

$$m = -0.95955, \text{ and } c = 1.53398$$

Therefore, by fixing the mean of the Normal distribution to zero and using the above formula the expected queue distribution of the play-out buffer can be obtained. By examining the tail of the distribution, one can determine the approximate queue size for a particular CLR. Figure 5.11 shows the optimum play-out queue delay for a CLR greater than $1 \cdot 10^{-8}$. These results relate the bit rate of the connection to a particular queue size. The queue size of the play-out buffer is smaller for lower bit-rate connections; however, due to the lower service rate of the buffer, the cells would be held for a longer time in the queue.

Therefore, to reconstruct the deterministic inter-arrival time, more delay has to be added to lower bit rate traffic, while larger play-out queues are needed for higher bit rate connections.

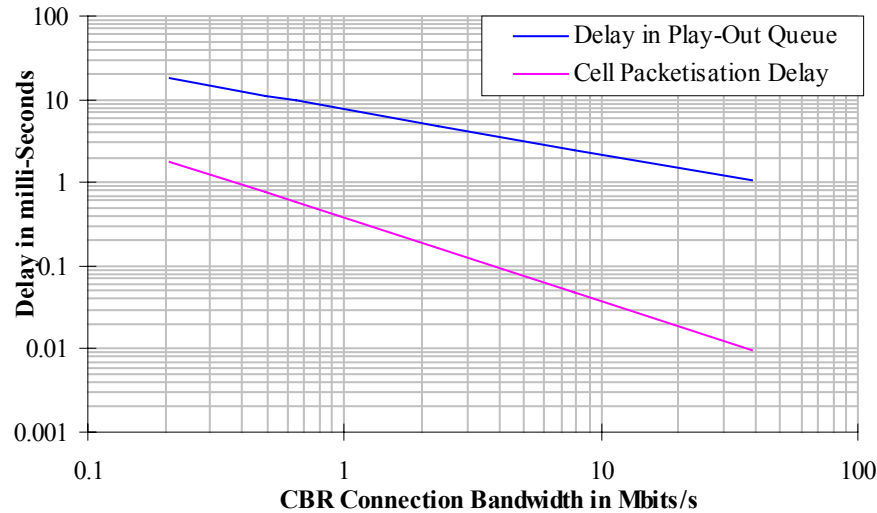


Figure 5.11: Delay within a CBR End-Terminal.

In addition to the delay due to play-out buffering, a cell will have packetisation and de-packetisation delay in the end-systems. AAL1 packetisation delay has been added to figure 5.11 to complete the total delay caused to a cell stream by the receiving end systems for CBR sources.

5.2.5 Summary

Simulations were undertaken to investigate the optimum buffer length for an arbitrary bit rate. The simulations determined that lower bit-rate traffic is less susceptible to CDV in terms of queue length. This is due to the fact that, lower bit-rate traffic has a greater packetisation and de-packetisation delay than higher bit-rate traffic and therefore needs less queueing space in the play-out buffer. Hence, by examining the arrival-pattern of “worst case” CBR traffic profiles it was possible to determine, by simulation, the play-out buffer dimensions. The arrival rate was then used to determine an analytical solution for the queue distribution of a worst case scenario. With the introduction of the play-out buffer additional delay was introduced. This delay reduces the QoS perceived by the user, particularly for interactive services. For example, when using a 210 kbit/s CBR application the delay added by the end terminals is 22.2 ms; this is before any delay due to the transmission system is included.

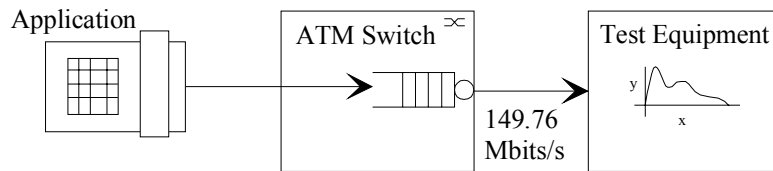
5.3 Experiment 3, Mapping Subjective QoS into Network Performance Parameters.

5.3.1 Introduction

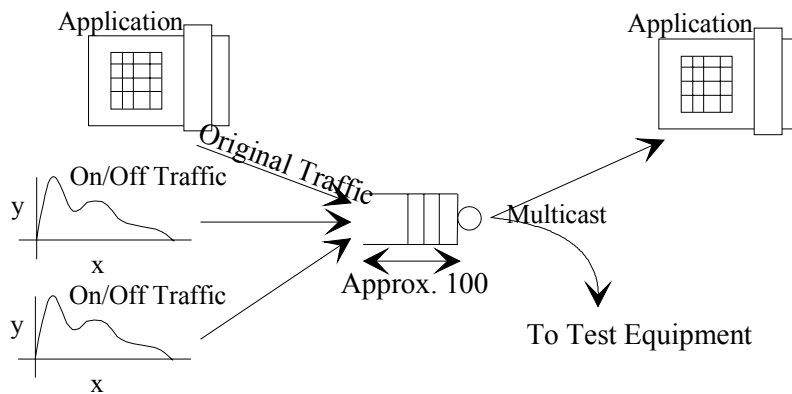
The objective of this experiment is to gain an assessment of QoS under different network conditions. Mapping QoS to network performance parameter will allow the reverse process of determining the QoS based on the network performance degradation in later experiments. In Chapter 4, it was stated that every end-system will be affected differently by CTD, CDV and CLR; some applications may have a click in the sound, other applications may need a complete system reboot. Therefore, the result of Experiment 3 is only valid for the applications considered; no experiment of this type can be applied to the general case. However, an assessment must be made between QoS and network performance in order to understand the effects of worsening network performance on higher layer QoS.

Four applications that need to transfer data over real-time circuits are tested to determine the maximum allowable degradation to the traffic characteristics, while maintaining the minimum acceptable QoS. First, each application was examined to determine the traffic characteristics at the input to the network; the functionality of the traffic sources was described and the resultant IAT graphs were plotted in figures 5.13, 5.14, 5.15 and 5.16. Second, the traffic from each source is tested under different conditions to determine the maximum allowable CLR and CDV before the QoS becomes unacceptable to users. The effects of worsening network performance are described in terms of a subjective assessment.

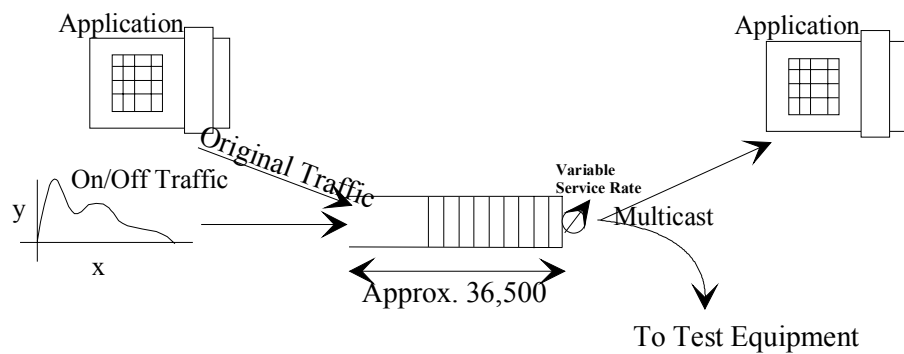
Experiment 3 had three stages. Firstly, the configuration of figure 5.12a was used to characterise the real ATM sources. Secondly, the traffic was multiplexed in a small buffer to determine the maximum CLR of the source, (figure 5.12b). Finally, the cell stream was multiplexed through a large buffer with background traffic to determine the maximum allowable CDV, (figure 5.12c).



a) Configuration to Determine the Traffic Characteristics.



b) Configuration to Determine the Maximum CLR.



c) Configuration to Determine the Maximum CDV

Figure 5.12: Experimental Set-up

5.3.2 Application Description

5.3.2.1 64 kbit/s N-ISDN Service

N-ISDN application uses a standard set of N-ISDN telephones, which are directly connected to an ATM terminal adapter. The application's data stream is packed into AAL1 PDUs, then transmitted on the ATM link from the terminal adapter (which has a line rate of 155.52 Mbit/s). The application uses 2*64 kbit/s

channels, each channel having a mean inter-arrival time of 2154.907 cells on a transmission link of 155.52 Mbit/s. The source was characterised on an unloaded network, having passed through one pure ATM link (155.52 Mbit/s), a single switch, and one STM1 (149.76 Mbit/s) link before entering the test equipment. Figure 5.13 shows the inter-arrival time density of the source. It is a CBR ATC cell stream with a standard deviation (σ) of 1.229248.

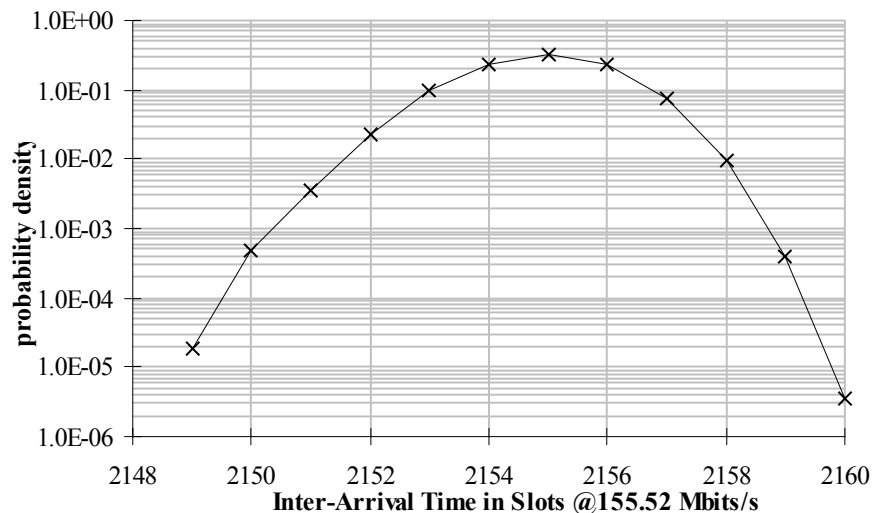


Figure 5.13: 64 kbit/s N-ISDN Telephony

5.3.2.2 384 kbit/s High Quality Audio Terminal

The High Quality (HQ) audio equipment has a similar hardware configuration to the N-ISDN terminal adapters. An audio signal enters an ISO (MPEG) D/S11172-3 codec and is fed into an ATM terminal adapter. This terminal adapter uses a simple FEC code, this being based on the ATM cell (AAL1) sequence number and a parity check. The FEC mechanism increases the bit rate slightly by 1/8th, but allows one cell in a block of eight to be corrected. The codec has a variable coding rate and thus can generate different CBR bandwidths: 64kbit/s, 128 kbit/s, 256 kbit/s and 384 kbit/s. In the experiments, the 384 kbit/s bandwidth setting at the codec was used. Once the terminal has generated the parity bits, the whole block of cells are transmitted using AAL1 adaptation layer functions. The inter-arrival time of the source is 314.2588 cells on a transmission link of 155.52 Mbit/s. The results are shown in figure 5.14. The traffic was analysed

through the same path as the N-ISDN and resulted in CBR ATC traffic stream with a standard deviation (σ) of 1.166809.

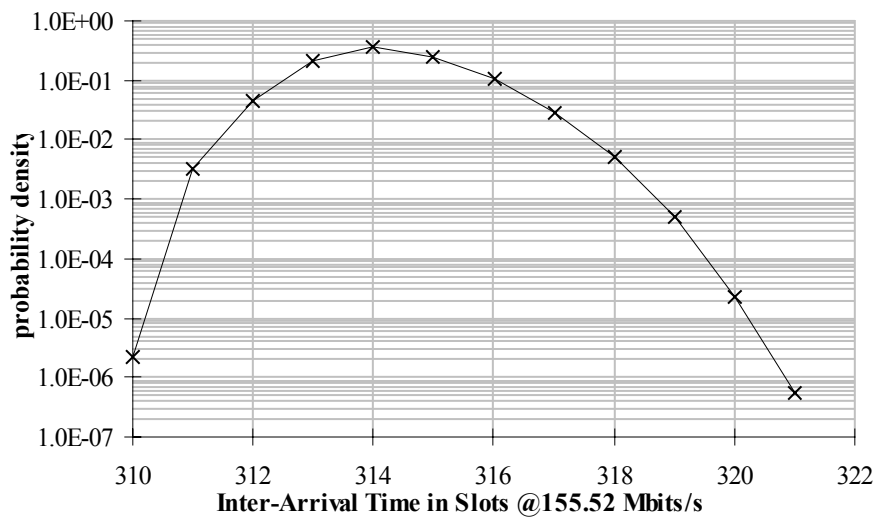


Figure 5.14: 384 kbit/s HQ Audio Equipment

5.3.2.3 H261 2 Mbit/s Video Conferencing

The 2.048 Mbit/s video conferencing was used to assess the QoS of a higher rate N-ISDN application subjectively. The H261 terminal was connected to the 2 Mbit/s PDH port of a broadband switch, this switch encapsulated the 2 Mbit/s stream into AAL1 PDUs and forwarded them onto the ATM network. This meant that the application did not generate cells, but an intermediate switch (located in the CPN) converted the incoming data stream into ATM cells. This configuration will be common with the introduction of ATM when existing protocols are encapsulated and transferred over an ATM network. The application was analysed at an SDH port of the switch and the inter-arrival time probability density is plotted in figure 5.15. The inter-arrival time of the application was 67.34064 cells on a transmission link of 155.52 Mbit/s and had a standard deviation (σ) of 1.827851.

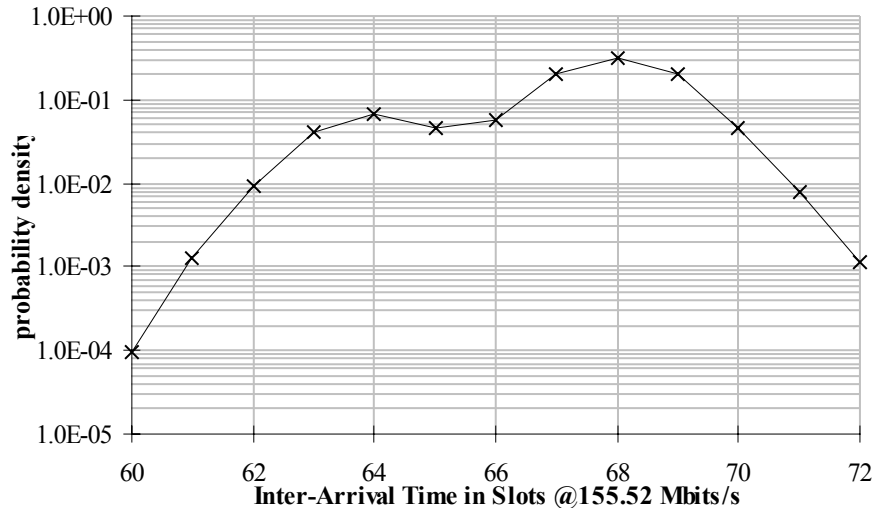


Figure 5.15: H261 2.048Mbit/s N-ISDN

5.3.2.4 1.716 Mbit/s ISABEL Multimedia Terminal.

The ISABEL application was designed for the interconnection of audiences. ISABEL was developed primarily for the RACE and ACTS summer schools, where several auditoriums had to be interconnected to create a large international distributed virtual auditorium. This enabled the attendees to achieve a sense of participating in an event independent of its location.

ISABEL is a Computer Supported Co-operative Work (CSCW) application that provides the basic technological framework for supporting remote collaboration in various areas of professional activity. ISABEL has been designed as a configurable CSCW environment that supports several interaction modes. The areas where ISABEL has been used successfully have included:

- Tele-education/training.
- Telework.
- Telemeeting.

The application is based on SUN workstations and uses UDP/IP over an ATM protocol stack. It has a variable bit rate characteristic and, having an IP protocol stack, utilises all of the allocated bandwidth in the ON state of the burst. The

mean rate of the application is 1.716 Mbit/s, but this bit rate can vary significantly with changes in screen size, audio quality, picture refresh, audio talk-spurts and picture content. The application was characterised using test equipment and the IAT of the Isabel application can be seen in figure 5.16. Within figure 5.16, two characteristics are presented, one when no shaping device is used, and one with the shaping device set to a rate of 2.2 Mbit/s. This feature will be highlighted in Experiment 4.

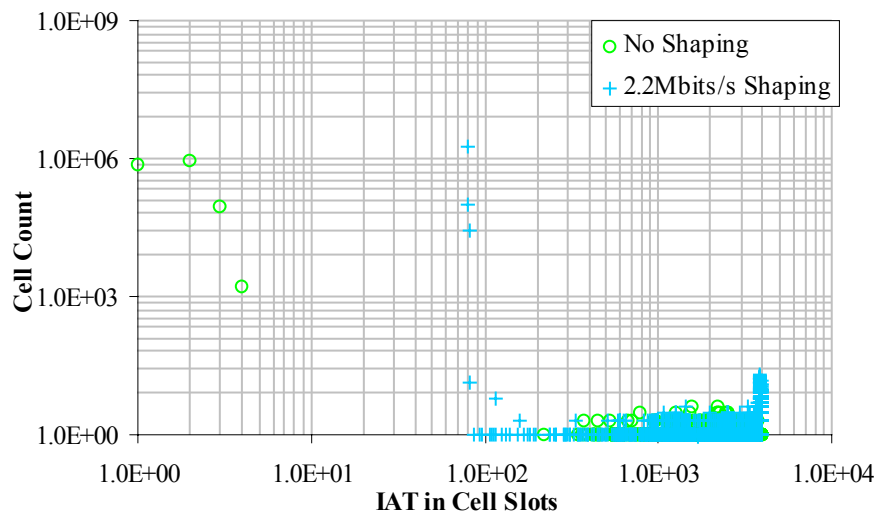


Figure 5.16: Isabel Multimedia terminal.

5.3.3 Results

5.3.3.1 Maximum Cell Loss Ratio

An assumption is made that the CLR in the buffers is the total number of cells lost divided by the cells transmitted. These measurements do not take into account SECBR, CMR or CER.

An estimate of all the application's CLR performance has been researched. This was done subjectively by listening or viewing an application and determining a measure based on the user's acceptance of the service. A number of people were elected to be "guinea pigs", each person using the application under decreasing QoS. When the QoS reached a point where the application still worked, but the errors were bordering on being irritable, network performance measures were

taken of the application's traffic. This produced a mapping between the minimum QoS and network performance. This mapping can be seen in table 5.6 where the minimum suitable QoS can be seen in terms of the CLR.

Application	Cell Loss Ratio
H261 Video Conferencing	$1.85 \cdot 10^{-5}$
384 kbit/s HQ Audio ⁴	$4.57 \cdot 10^{-5}$
64 kbit/s Audio	$3.42 \cdot 10^{-4}$
Isabel Multimedia Terminal	$1 \cdot 10^{-4}$

Table 5.6: Minimum QoS to CLR Mapping

CLR is a difficult dimension to give to any service class, as the application's QoS depends on many factors with regards to losses: the distribution of cells lost within a connection; the types of end system; and whether, a statistically significant number of cells have been transmitted to obtain an accurate measurement. The last point is important for experiments: to obtain the $1 \cdot 10^{-4}$ CLR for a single telephony circuit would take approximately 2 hours. The mean duration of a telephone call is between 2-3 minutes, so, it would be quite difficult to measure small values of CLR on real network connections, although small values of CLR are assumed.

The audio devices (HQ Audio and N-ISDN) suffered a range of imperfections (from audible clicks to loss of sound for small periods) with cell loss. The losses in sound were assessed to be far more important for the high quality audio application than the telephony application. This is due to the *expectation* the user had for that particular service. The users were comparing the high quality audio application to a CD player, hence, the users expected to hear CD quality sound and therefore a higher quality comparison was made. The telephony circuit was examined and the users found that a lower QoS was acceptable, as the information content in the conversation was more important than the audio quality, a conversation was still understandable with 8% of uniformly lost cells. In addition,

⁴ The audio application had an FEC and so a single error appearing in a block of eight was corrected.

the comparison between existing POTS and the telephony application meant that the users expected a low quality application and were quite resilient to clicks in the audio component. The telephony circuit was assessed for voice only; data and/or modem requirements have not been examined, but might be expected to be more stringent.

The H261 2 Mbit/s video conferencing unit was assessed for both audio and video quality. Due to the method of implementation, a cell loss would effect the audio and video components of the application at the same time. When assessing the quality, the “panel” found that the quality of the audio had an equal if not greater importance than the quality of the video. This demonstrates that the lower bandwidth audio application contains as much information to the “human” user as the higher bandwidth video component.

The ISABEL terminal is a real-time VBR application that uses UDP/IP. A subjective assessment determined that picture flicker was not as irritable to the user as the clicks and losses in sound. It seemed that the audio component was, again, far more important than the video one. A cell loss ratio of $1 \cdot 10^{-4}$ gave an acceptable level of QoS. However, the parameter of CLR is a very rough measurement as the effect of the CLR is dependent on the distribution of the cell losses. If the IP layer does not receive a complete packet, the data for that packet is discarded, i.e. if one cell is lost within an IP packet then the whole packet is discarded. The same single packet discard would result if the distribution of cell losses were clumped and several cells were lost in one IP packet. The disturbance to the application is the same, i.e. one packet is lost, despite the different CLR's.

In table 5.6, the values of cell loss allows an estimation of the human user's requirements from the network. When the end user is a computer, the requirement for low CLR will become more significant. The network operator is not aware of which application is being used across the network, so the operator must assume the worst case and dimension for the most stringent communication requirements.

Over a network, all real-time connections must take on the most stringent requirements of performance. The specification of stringent performance

requirements are difficult to determine as they are based on: the bit rate of the connection; the requirements for re-transmission; sensitivity of the higher layers to cell losses; the distribution of cell losses; and finally, manufacturer dependent features. As each application is affected in a different manner, the operator must assume that all traffic needs a high performance circuit. Only the user can take on the responsibility of making gains through statistical multiplexing of real-time traffic, as increasing the admissible load boundary usually causes higher CLR's. This should only be allowed in the CPN's, where the users have direct control over both the end-systems and switching elements.

5.3.3.2 Effects of CDV

Application	peak-to-peak CDV
H261 Video Conferencing	19 cells
384 kbit/s HQ Audio	11 cells
64 kbit/s Audio	2 cells
Isabel Multimedia Terminal ⁵	111 cells

Table 5.7: Minimum QoS to CDV Mapping

When assessing an audio application's QoS, loss of sound becomes the significant effect due to CDV performance degradation. In table 5.7, the maximum peak-to-peak CDV can be seen for the applications under test. The perceived QoS caused by CDV was indistinguishable from the effects due to a high CLR. The reason for the loss of sound due to CDV is that the play-out buffer (see Experiment 2), will reset or flush depending on the deviation of the cell stream. Thus, when the buffer is either in the empty or in the overflow state, cells would fail to arrive on time and a gap would appear in the information stream to the higher layer. Similarly, losses of picture and/or sound happened in the video conferencing application

⁵ 4 Mbit/s shaping rate was used to introduce CDV into the cell stream. Isabel is based on IP and so the whole packet needs storing before the data continues to the next higher layer. This means that the last cell in the IP packet has the greatest CDV, see Experiment 4 for further details.

when too much CDV occurred. Again, the effects observed are similar to those caused by cell loss.

When using these applications, the effects of too much CDV and cell loss are the same at the application layer. The application generally had faults directly related to the amount of CDV or cell loss, until a point where a slight increase in CDV or CLR caused the application to fail completely. This resulted in a “cliff edge” type failure, the application working with an acceptable QoS at one instant and then crashing when a small decline in network performance occurred.

Adding CDV to the ISABEL multimedia terminal with the introduction of a shaper is described in Experiment 4. The subjective QoS assessment of this application with increasing CDV found that the picture slowed, as picture frames were lost and again, gaps appeared in the audio component. This application was quite resilient to CDV and heavy shaping could be employed before the application failed.

5.3.4 Summary

In this section, Stringent Class 1 applications were characterised at the ATM layer to determine the bandwidth and cell variation around a nominal inter-arrival time. The end-systems were described to facilitate an appreciation of the types of application investigated.

The applications were assessed for QoS degradation with increasing CLR and CDV. When cell loss is present, an application can still be acceptable to a human user. Different services have different expectations: for example, clicks during a CD quality audio transmission are more annoying to the user than clicks during a telephony connection. Applications can afford to have a high CLR when the end users are human; but, when the end user is a computer then slight data losses may be significant. As the CDV increases, cells are effectively lost from the “empty” and “overflow” states of the play-out queue. Thus, to increase the performance of circuit emulation equipment with respect to the CDV, longer play-out queues are needed at the input of the end-system. However, increasing the queue size also increases the effective delay across the network, as stated earlier in Experiment 2.

5.4 Experiment 4, PCR Reduction by Traffic Shaping while Maintaining QoS

5.4.1 Introduction

5.4.1.1 rt-VBR Source Model

As stated in chapter 4 there are no native ATM applications on the market and there will probably be none in the short to medium term time scales, [SCHM97]. Therefore, this thesis uses IP applications over ATM for experiments on the concepts of shaping traffic.

Data is written to the IP protocol stack by higher layer protocols until no data remains or the maximum MTU size is reached, whichever occurs first. The IP layer will forward the data block with the IP header to the ATM network interface card (NIC). Once the NIC receives this data it will operate as quickly as possible to attach an AAL5 trailer, segment the AAL frame into cells and send the fragmented data onto the physical layer with the appropriate ATM header.

Figure 5.17 shows an illustration of an ATM traffic trace from a NIC that uses a TCP/IP protocol stack.

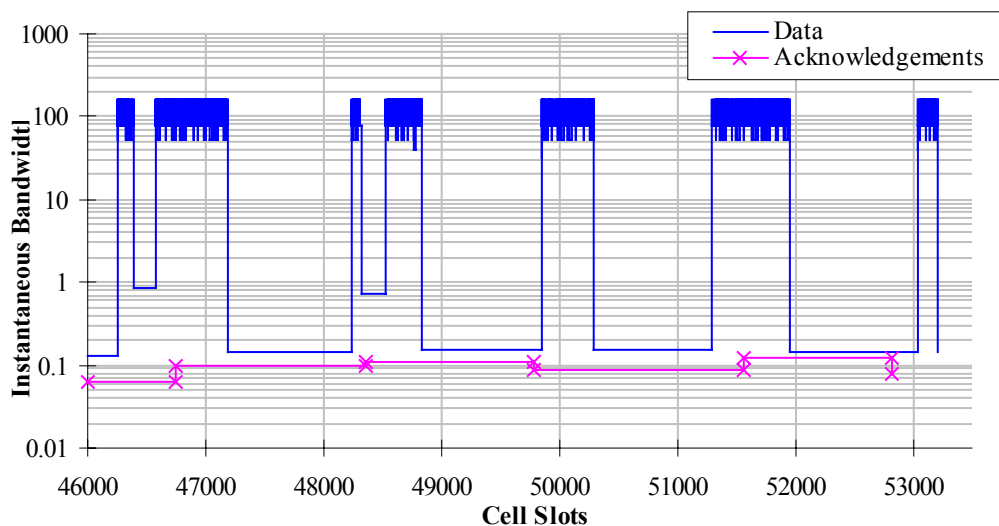


Figure 5.17: IP data traffic, Peak 155.52 Mbit/s, Mean 22 Mbit/s. (SDH TX Line.)

The IP over ATM process has a large packetisation delay because of the IP layer. Therefore, the traffic originating from ATM NICs is very bursty in nature. This

leads to non-optimal multiplexing solutions when control is not present, particularly over rt-VBR ATCs. Similarly, over-controlling a traffic source to increase the network's multiplexing density can reduce the performance of an ATM network to a level where a cheaper solution, e.g. 100baseT Ethernet, could provide the required service.

5.4.1.2 Concepts of QoS, Performance and Shaping

Due to the random nature of many applications, traffic contracts are difficult to estimate and rules-of-thumb are currently used to order ATM circuits. The main specifier is currently the PCR, (the SCR and MBS can be difficult for users to specify as a more detailed knowledge of the traffic characteristics is required). Introducing a simple PCR traffic shaper maintains this simplicity and keeps the traffic within the negotiated bounds.

When a connection is ordered, the user preferably selects the peak rate of the application, as this will allow the cells to be transferred across the network with the smallest delay. For an IP application that corresponds to the peak rate of the NIC, perhaps 155.52 Mbit/s, as seen in figure 5.17. This may be preferable speed-wise, however is not at all economical. Thus, the user will then re-evaluate the network connection and apply a QoS/Price limit. By reducing the peak rate and applying shaping the user can still obtain acceptable QoS at reduced communications costs.

To highlight the need to optimise VBR traffic to reduce communications costs, a conceptual graph has been produced (figure 5.18). On this graph, the dotted line Y1 represents the instantaneous bandwidth usage. The graph shows that the peak rate necessary for the application is 31 Mbit/s. In addition, one can judge, for example, 26.5 Mbit/s is required to transmit all but 10% of the application's traffic, without shaping. By applying shaping, the peak rate requirements can be lowered to an acceptable QoS point, which lies just above the QoS "cliff edge", (line Y2).

This method enables a user to reduce the application's bandwidth requirement to while maintaining an acceptable level of QoS. As the operator guarantees only

network performance and not QoS, the user can determine a suitable QoS for the application based upon saving made in communication costs.

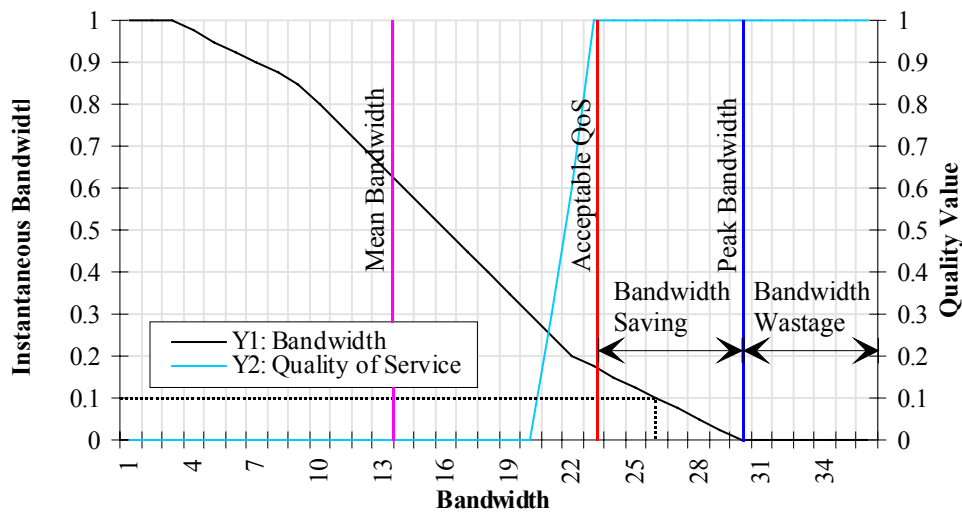


Figure 5.18: VBR Optimisation Concept Graph

In figure 5.18, Y2 is the QoS grade, the value 1 representing excellent QoS, and 0 very bad QoS.

Applications fail abruptly when the gradient between excellent QoS and very bad QoS is steep. Thus, the user must know the minimum bandwidth that the application requires: this is at the QoS “cliff edge”. If the user specifies a bandwidth below this edge then a connection charge will be generated without the application working. For example, in figure 5.18 the choice of 15 Mbit/s would generate a connection charge despite the fact the application will not work. The QoS “cliff edge” is normally determined through trial and error.

There is also a requirement on the operator to control the network very accurately once a user has determined the optimal point for the peak-rate bandwidth.

Degradation in any of the network performance parameters could take the application over the QoS “cliff edge” and cause the application to fail. Thus, once an operator has accepted to guarantee a traffic contract, the network performance must be maintained throughout the lifetime of the connection. The only way to improve the QoS of the application, due to worsening network performance, is for

the user to order more bandwidth. As the user increases the bandwidth, CTD and CDV degradation caused by the shaper will reduce, hence offsetting the degradation caused by the network. From Experiment 1, it can be seen that using international and intercontinental CBR channels, with a peak-rate allocation policy, the traffic characteristics change very little across the ATM network. This indicates that once the CBR channel's traffic contract has been accepted, the deviation from the contracted parameters will be very small, at least with the network loaded as it was during these experiments.

5.4.2 Experiment Set-up

5.4.2.1 Configuration of the FTP Trial

A simple data transfer mechanism (FTP) was tested using the proposed VBR optimisation, which is based on single rate shaping. FTP runs on top of TCP/IP and requires an asymmetric bi-directional ATM connection. A 5 Mbyte data block was transferred repeatedly between two workstations (Sun Sparc 20) and in each trial, the output rate of the traffic shaper was adjusted. In figure 5.19, the configuration of the FTP experiment is shown. The data traffic has the highest throughput of information and so this channel is shaped; the acknowledgement channel had a much lower bandwidth and so shaping was considered unnecessary during the experiment.

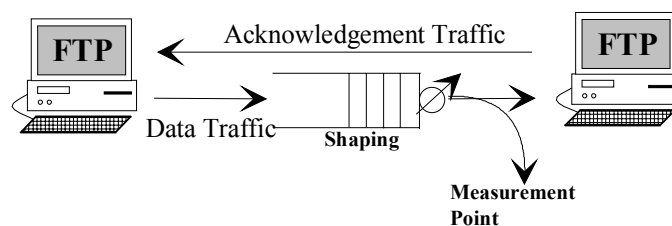


Figure 5.19: FTP Optimisation Experimental Set-up

It is not easy to measure the QoS of most applications quantitatively. However, an FTP session can be assessed. The QoS parameter defined for an FTP application is the time to transfer the 5 Mbytes of data from one workstation to another.

Within Experiment 4, the file was transferred from one workstation to another and

the transfer time was recorded. The shaping rate was reduced and the experiment was repeated. The results are presented in section 5.4.3.1

5.4.2.2 Configuration of a Multimedia Terminal Trial

The ISABEL Multimedia Terminals, [ISAB97], have proved themselves very popular for trials and demonstrations and their traffic has been analysed extensively in research, [DEL15] and [BAUM98]. The terminals support both point-to-point and point-to-multi-point connections. The application runs on Sun workstations and uses a UDP/IP protocol stack to transfer data. TCP/IP is not used, as it does not support IP multicast connections.

The application consists of audio, video, scratch pad, pointer and slide presentation devices. The audio can be configured to provide a number of different qualities, from low bit-rate sound (8 kbit/s) to high quality sound (600 kbit/s). The audio component uses silence suppression to reduce the volume of traffic when the microphones are not receiving audible noise. The video uses moving JPEG to encode the pictures and hence there is no predictive coding of the video stream. The number of frames per second and the size of the video window can be altered; both of these adjustments have a marked effect on the application's mean bit rate. The scratch pad and pointer have little effect on the traffic characteristics and were not considered in the experiments. The slide presentation function is more complicated as the slide presentation is transferred using FTP before the presentation, and the speaker then uses a control mechanism to sequence the slides at each end.

A standard configuration was defined use in the sets of experiments: a large screen size with a refresh rate of 12 frames per second and a high audio quality of 20 frames per second. From this, a mean rate of 1.716 Mbit/s was transmitted from the application.

A multimedia terminal is considered to be a real-time VBR source, as the data needs to be transferred as quickly as possible. The application has been tested over local, international and intercontinental ATM links and shaping has been

used to obtain the minimum operating bounds, in an effort to minimise the ATM bandwidth required.

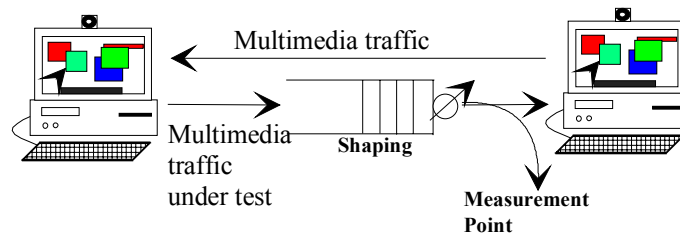


Figure 5.20: Multimedia Optimisation Experimental Set-up

ISABEL multimedia terminals were set-up as in figure 5.20 and one link was tested to determine the effects of shaping on the application. QoS assessment was more difficult to test on this application than with the FTP as no objective assessment could be made, so a panel of subjective testers was used. When the applications became noticeably delayed, or the audio and video contained errors, the application was declared unacceptable and the bandwidth value of the shaper recorded. From the recorded traffic, the effect of the shaper on the traffic characteristics is observed.

5.4.3 Results

5.4.3.1 Optimisation of a File Transfer Protocol

From figure 5.21, it can be seen that the best possible transfer time (and hence QoS) is 3.6 seconds. To achieve this transfer time it is immaterial whether a peak bit rate of 25 Mbit/s or 155.52 Mbit/s is used. This means that the user can make substantial savings in communication costs by reducing the peak rate of the connection, as there is no difference in QoS between the bandwidth values of 25 Mbit/s to 155.52 Mbit/s.

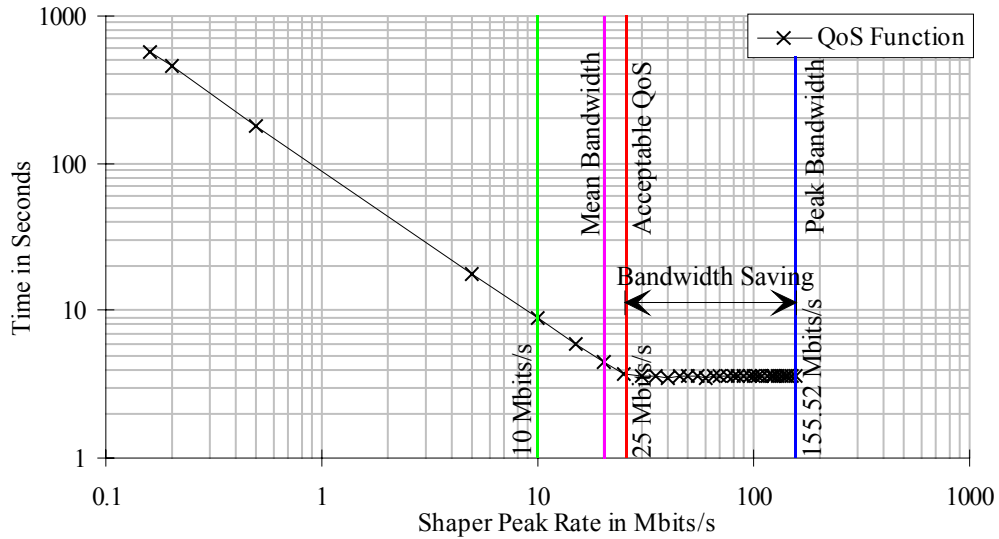


Figure 5.21: Time taken to Transmit a 5 Mbyte Binary Data Block.

The bandwidth was reduced even further to observe the effects of shaping on the QoS. The transfer speed decreased as shown in figure 5.21. The mean rate of the unshaped data traffic was 20 Mbit/s and when the shaping rate reached this value cell loss was expected because of buffer overflow in the shaper. However, as TCP/IP is adaptive to the congestion of a network, the traffic shaping caused the terminals to reduce the transmission rate, causing no cells to be lost. The time to transfer the 5 Mbyte data block increased inversely proportional to the available bandwidth.

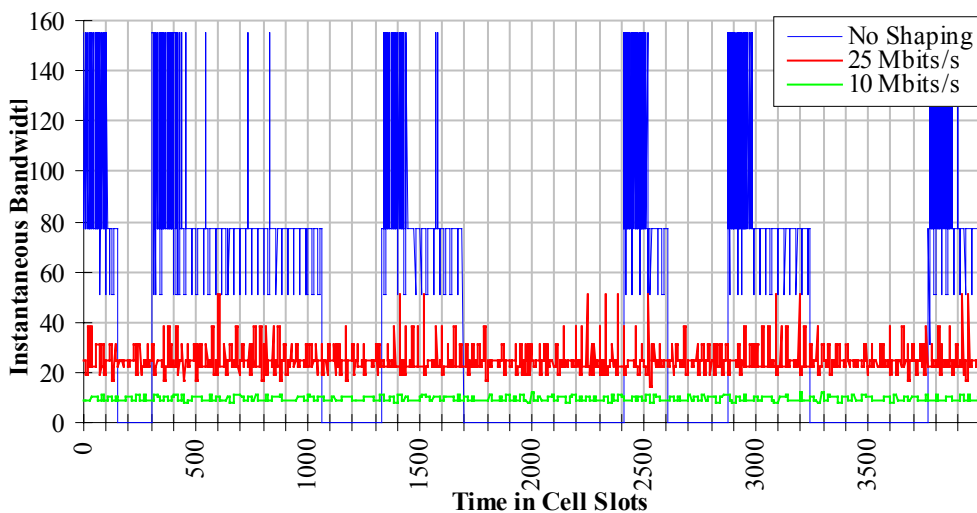


Figure 5.22: Shaping an FTP session (Data channel only).

The FTP traffic on the ATM network was observed using test equipment and the data traffic is shown in figure 5.22. From the diagram, the unshaped traffic using all the 155.52 Mbit/s link capacity appears very bursty. Reducing the connection bandwidth to 25 Mbit/s alters the traffic characteristics quite substantially. It can be seen that shaping will lead to an increase in the traffic predictability, which may in turn lead to more accurate CAC and enable a higher multiplexing density across the network. When even more traffic shaping is applied, i.e. a bandwidth of 10 Mbit/s, the application continues to transfer the data, but at a reduced rate. The application uses all the allotted 10 Mbit/s set by the shaper and so, offers both the user and network operator a method of independently controlling TCP/IP traffic within the customer premises network or public network respectively.

5.4.3.2 Optimisation of a Multimedia Terminal

Shaping was applied to the ISABEL multimedia terminal (with standard settings) to obtain the minimum bandwidth and hence the optimal amount of resources required to run the application. Figure 5.23 shows the instantaneous traffic pattern on the ISABEL multimedia terminal. Similar traffic patterns were observed as in the trials with FTP; this is because the multimedia terminal used an IP over ATM protocol stack. The application could peak to 155.52 Mbit/s unshaped, but the mean bit rate was 1.716 Mbit/s.

Shaping could be heavily applied to the application and 2.2 Mbit/s was judged, subjectively, to be the minimum amount of resources needed to achieve a good QoS. Shaping the multimedia traffic below this value would cause the application to fail completely. Thus, the optimum bit rate was just above the mean rate, and at this rate the traffic from the shaper had almost a CBR characteristic.

UDP does not let the transmitter adjust its bit rate according to network congestion like TCP/IP; hence, the mean bit rate does not adapt. When a shaper queue of 36,000 cells is used at low service rates, the end-to-end delay becomes very large. Delays of 6 seconds were recorded between end-systems and subjectively this gave an unacceptable QoS. When the application was shaped too heavily, the QoS of the application quickly worsened. The picture flickered and audible clicks were

present in the sound. Applying subjective assessment, picture flicker was not as irritable to the user as the clicks and loss of sound and again it seemed that the audio component is far more important than the video one, as had been previously noted in Experiment 3.

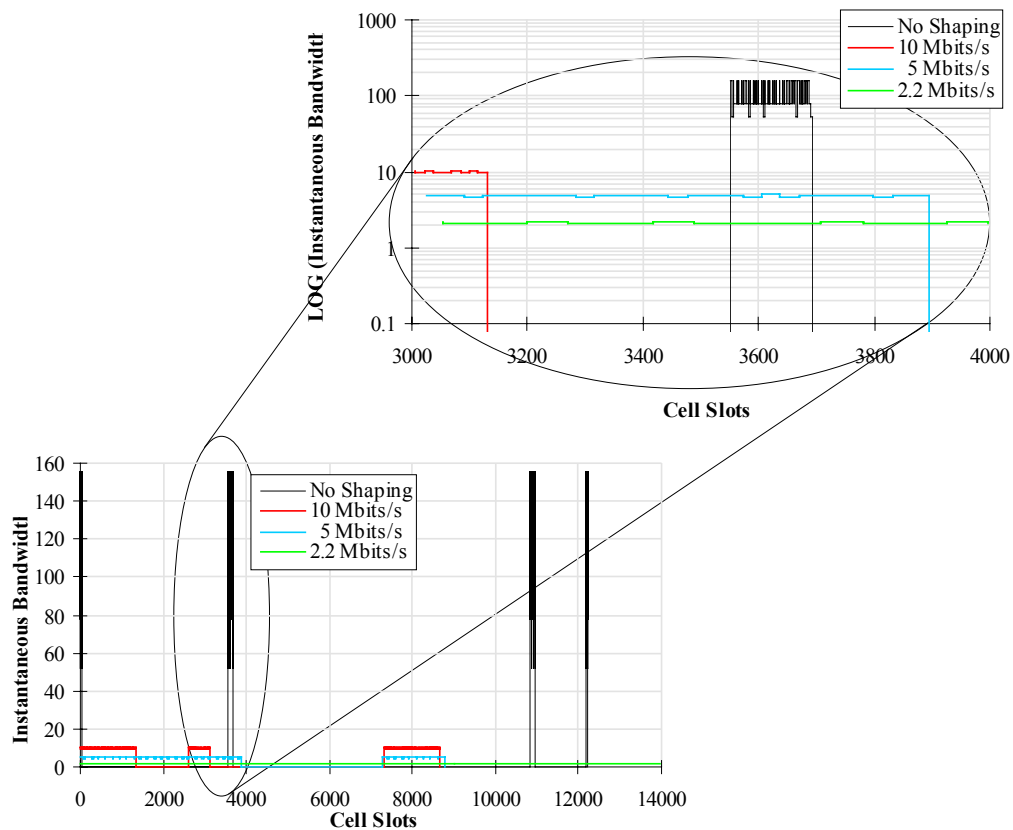


Figure 5.23: Shaping an ISABEL Multimedia terminal

When introducing a shaping device into the ISABEL terminal equipment, an assessment of the shaper maximum queue size is required. The benefits of reducing the burstiness of a connection should not be out-weighted by the expense of the memory and the delay caused by queueing. Using the model of the ISABEL terminal validated in [BAUM98], a simulation using the YATS simulator was set-up by the author, to investigate the length of the buffer necessary to shape the application. The traffic from the unshaped multimedia model entered a shaping device, the service rate of the shaper was varied, and the queue length was assessed. The results are shown in figure 5.24.

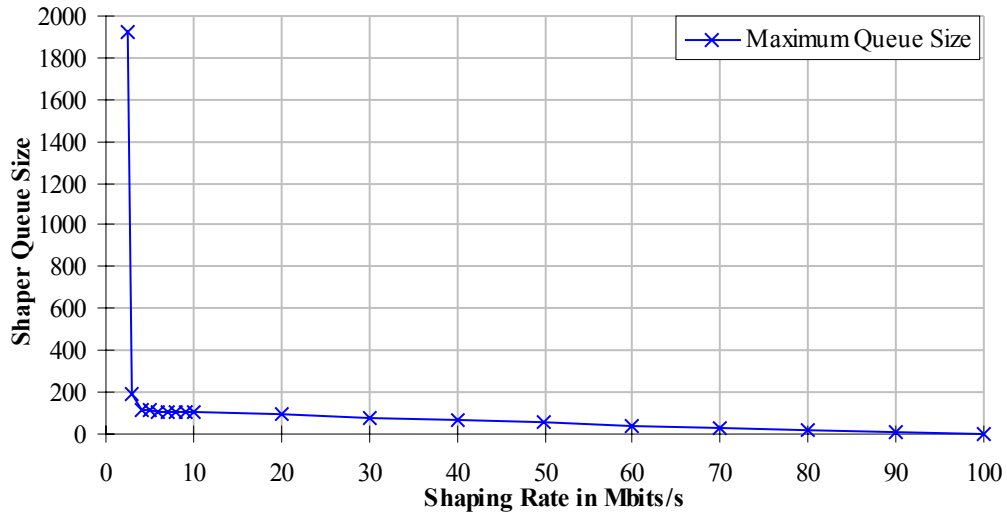


Figure 5.24: Maximum Shaper Queue Length

At a shaping rate of 100 Mbit/s, a queue of one is sufficient to transmit the application’s traffic without loss. As shaping is applied, the maximum queue size increases linearly until 4 Mbit/s where the queue size is 111 cells. When the shaping rate decreases below 4 Mbit/s the queue length required to shape the traffic quickly increases. At 3 Mbit/s the queue is 192 cells and at 2.5 Mbit/s the queue length is 1935 cells. The application has a mean rate of 1.716 Mbit/s. The reason for the large increase in queue size, is that, as the peak rate shaping reduces toward the mean of the application, it becomes necessary to fill every empty slot from the shaper to transmit all of the application’s data. However, as the application is very bursty in nature, some of the cells leaving the shaping device are empty. Thus, the number of cells necessary to transmit the mean amount is increased due to the burstiness. The minimum shaping rate of the application is selected to be between the mean rate and the peak rate of the application; this is called the sustainable cell rate and, from figure 5.24, is nominated to be 4 Mbit/s for the ISABEL terminal.

Having a shaper in the edge device increases the apparent delay across the network for the application. From figure 5.24, the ISABEL application has an optimum traffic shaper rate of 4 Mbit/s and at this rate the maximum delay caused by the shaping device is 11.807 ms.

The application uses IP over AAL5. This means the AAL5 frame needs to be transmitted as quickly as possible, as described in section 5.4.1.1. The shaper increases the inter-arrival time of cells, hence progressively increasing the CTD of each cell belonging to the AAL5 frame. The last cell of the AAL5 frame has the greatest delay. At the receiver, all the ATM cells belonging to an AAL5 frame are concatenated before being forwarded to the application layer. Thus, it is necessary to wait for the last cell of the frame, which has been delayed the most by the shaper.

Consider a burst of cells from an AAL5 stream that has n cells per frame. Each cell with an inter-arrival time of m cells and a shaper's service rate set to m' cells then,

- If $m \geq m'$ there is no additional delay caused by the shaper.
- If $m < m'$ the increase in delay (Δd) caused by the shaper is:

$$\Delta d = (n - 1)(m' - m) \text{ Cell Slots}$$

- If the Inter-Frame Time (I) of the AAL5 frames is small, work from previous frames will still be in the queue on the arrival of the frame. Thus, the additional delay caused by the shaper on the N^{th} burst is:

$$\Delta d = (m' - m)(Nn - 1) - (N - 1)I \text{ Cell Slots}$$

The ISABEL multimedia terminal is a dual process model, i.e. an audio component and a video component. The mean inter-frame time of the video component (I_1) is 3742 cells and that for audio (I_2) is 17660 cells. Thus, at a shaping rate of 4 Mbit/s (i.e. $m' = 38.88$ cells) the worst case delay is when residual traffic is left from the arrival of an audio frame (35 cells) and video frame (87 cells) in the queue and a new video frame (87 cells) arrives. From the analysis a delay of 11.849 ms is calculated, which approximately matches the delay observed in the simulation of 11.807 ms. The re-packetisation process needs to wait a further 11.807 ms for the reception of the last cell in the packet. Then the total delay caused by the shaper to the cell stream is 23.614 ms. Further analysis

on the performance of a shaping device with data connections can be found in [NIES93].

This delay of 23.614 ms would be equivalent to having an additional 8660 cell spaces (@155.52 Mbit/s), distributed within queues of switches along the path of the call. This may mean that the application could be better suited to the nrt-VBR category, although buffer sizes of 115,000 cells are expected per switch for the nrt-VBR ATC, [BRIE98]. More buffering is not recommended along the real-time path, as synchronous applications are very dependent on the CDV across the network. However, if the buffer has service rate of 4 Mbit/s, it only needs a queue of 111 cells, making it a cheaper device to manufacture. In addition, the shaper is the responsibility of the customer and allows the customer to vary the delay in accordance with the QoS required.

When viewing the terminal and applying shaping, 2.2 Mbit/s was subjectively deemed to be optimal for the QoS/bandwidth trade off. However, when the delay was assessed through simulation, it was found that very large queues are required to buffer this low shaper rate. Hence, a better value based on the analysis of queue lengths is of 4 Mbit/s, which had the optimal combination of delay and bandwidth. The selection of either 2.2 Mbit/s or 4 Mbit/s will depend on the cost of the bandwidth, the cost of buffering, and the irritation caused by delay between the two end terminals. If the end-systems are close (low network delay), then a larger delay in the shaper can be accommodated making 2.2 Mbit/s suitable. If the delay over the network is large then increasing the bandwidth to 4 Mbit/s maybe a better solution, as an interactive conversation would be easier with reduced delay.

5.4.4 Summary

Optimisation of VBR sources using a shaper is important. The user and service provider can reduce their costs by applying shaping to a stream of cells, from which savings can be made in the resources required while maintaining a good quality of service. In previous work, traffic characteristics have assumed to be fixed with the network having to adapt to the application's profile. Modelling the traffic profile for different applications has been a topic of much research (for

example MPEG [AARS97], Multimedia terminals [BAUM98] and WWW traffic [FARB97]) but, maintaining the traffic characteristics may prove in most cases to be unimportant, [DEL10], [DEL15] and [GRAF97]. Within this thesis, it has been shown that traffic characteristics can be moulded into cheaper CBR ATCs without loss of QoS.

The most effective method of control available to the operator is the charging mechanism. Once the mechanism is implemented, users can quickly adapt their traffic (traffic shaping) to obtain the largest possible savings. However, charging mechanisms cannot be complex, and simplicity and predictability are important factors, [RAFF96]. Thus, users can be influenced by economic pressure to incorporate a shaping device in the transmitting end-system.

The network operator should only guarantee network performance parameters across the network. These objective measures can describe, define and measure a connection between end points. The user can select different ATC's, and can mould the applications traffic characteristics into cheaper traffic descriptors, appropriate to the required QoS. Therefore, QoS should then remain the user's responsibility and network performance the responsibility of the operator.

5.5 Experiment 5, Increasing Multiplexing Density by Traffic Shaping

5.5.1 Introduction

VBR sources are quite difficult to handle over an ATM network, particularly over real-time circuits as small switch buffers have difficulty absorbing burst level queueing. Recommendations exist for network performance parameters over ATM networks, which allow the specification of minimum and maximum values for CTD, CDV and CLR on the ATM layer. Initial network dimensioning will allow the operator to guarantee the fixed parameters; however, CDV and CLR are dynamic parameters that change with the network load. The network operator is able to control these performance parameters with Connection Admission Control. CAC is used to control the admission of new connections onto the network and thus gives the operator some control over the dynamic, time varying parameters.

From ITU-T Recommendations, the maximum delay that is expected through an ATM switch is 300 μ s. An ATM switch has both processing delay and queueing delay, and from experiment 1, delay measurements show that buffers in each switch were of the order of 76 cells.

Applications like multimedia terminals and video-on-demand will probably be the first real-time network services. Many of these applications are CBR in nature and can be readily admitted onto the network. CAC mechanisms are well understood for CBR services, which cause cell scale queueing, for example [DOMI94]. Predictions based on the M/D/1, ND/D/1 and Σ ND_i/D/1 models can be used to provision for suitable CDV and CLR guarantees [DEL30]. One multimedia terminal (ISABEL) is used within this experiment to investigate multiplexing of real-time VBR services. The traffic characteristics of the ISABEL are well understood [DEL06] and have been modelled [DEL10] from the moving JPEG frame layer to the ATM layer to produce accurate, memory orientated models. This model has an ON/OFF traffic characteristic. This means that obtaining any effective burst scale queueing gain (or statistical multiplexing gain) with small buffers from ON/OFF traffic is limited.

Optimising VBR sources for QoS and resource saving within an end-system may help to increase the utilisation of an ATM network by reducing the need to absorb burst scale queueing, [WINS97]. Simulations were constructed using the ISABEL multimedia terminal model, [BAU96a], to discover how many terminals are allowed on to an ATM network of 149.76 Mbit/s (STM1 transmission) and to assess the utilisation of the network. Queue sizes of 50 and 100 were used to represent queue sizes of real-time circuits.

5.5.2 Experiment Set-up

5.5.2.1 Application Description

There were several noteworthy features about the traffic from the ISABEL application. Firstly, the application is based on IP and uses an AAL5 protocol stack. When the AAL5 packet is filled, the whole frame is presented to the ATM

layer. The terminal had fixed AAL5 frame sizes, 4176 bytes for the video stream and 1680 bytes for the audio stream. The NIC transmits the data as quickly as possible utilising the maximum peak rate allocated by the shaper. Secondly, the mean rate varied according to the user's preferences on the terminal: large/small windows, low/high quality audio or high/low picture refresh rates. Throughout these sets of experiments a so called, "standard scenario" was used. This was a "large" screen size with a refresh rate of 12 frames per second and a "high" audio quality of 20 frames per second. From this, a mean rate of 1.716 Mbit/s was transmitted from the application.

5.5.2.2 Multiplexing real-time Multimedia Traffic.

As discussed earlier, multimedia traffic is likely to be a pioneering broadband real-time application. Experiment 5 simulates the network utilisation for a suite of multimedia terminals based on the model of the ISABEL multimedia terminal that was verified in [BAUM98]. Experiment 4 showed that shaping significantly reduced the peak-rate bandwidth of rt-VBR applications by allowing them to tend towards CBR traffic characteristics. In Experiment 3, the ISABEL application was shown to need a CLR less than 1×10^{-4} to give an acceptable QoS to the user. This value was used in Experiment 5 to determine the multiplexing density of a suite of homogeneous multimedia terminals. The simulated multimedia terminals were set-up as in figure 5.25. A single shaping rate was applied to each end-system and the maximum number of sources was determined for a particular CLR in the switch queue. Trials were repeated with a different shaper rate to construct the admissible load boundary (ALB). Multiplexing density has been defined as the maximum number of sources on the network.

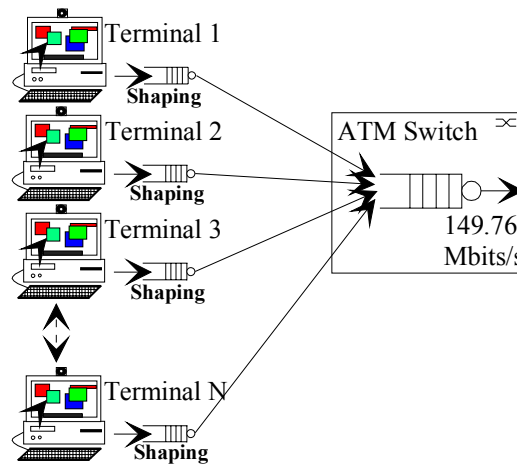


Figure 5.25: Simulation Set-up to Determine the Multiplexing Density

Using the ISABEL multimedia model, [BAU96b] and validated in [BAUM98], simulations were set-up by the author to determine network utilisation and multiplexing density. Using the results of Experiment 3, the maximum amount of shaping was found to be 2.2 Mbit/s (although this gave long shaper queue lengths) and $1 \cdot 10^{-4}$ was the maximum CLR for suitable end-to-end QoS. The shaping rates of the end-systems were varied from 100 Mbit/s to 2.2 Mbit/s and the number of sources were increased in each trial until a CLR of $1 \cdot 10^{-4}$ was reached within the multiplexer.

5.5.3 Results

Figure 5.26 demonstrates the number of ISABEL terminals that can be allowed onto a 149.76 Mbit/s real-time ATM network. Firstly, using an end-system shaper rate of 100 Mbit/s, with a multiplexer buffer length of 100, two ISABEL terminals are allowed onto the network before cell loss caused the application to have unacceptable QoS. Reducing the peak rate of the shaping device, and hence the burstiness of the source, allows the multiplexing density to increase. Using the optimum amount of shaping (4 Mbit/s as determined in Experiment 4) 38 multimedia terminals will be allowed to transmit simultaneously. Reducing the rate still further allows 51 sources (at 3 Mbit/s) to be transmitted.

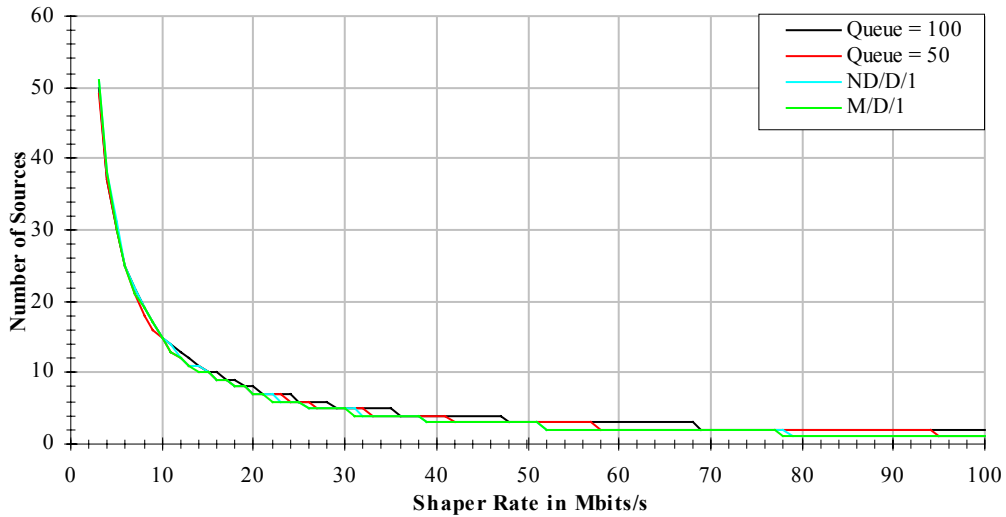


Figure 5.26: ALB for ISABEL MM terminal compared with ALB for CBR traffic

As the CDV tolerance of real-time circuits is limited, little statistical multiplexing can be gained. A comparison was made with standard CAC mechanisms for CBR sources. M/D/1 and the ND/D/1 analysis was applied to the VBR application using the peak rate of the shaping device as the traffic capacity of an equivalent CBR source. From figure 5.26, a close match is observed between the CBR analysis and shaped rt-VBR results. Thus, in terms of effective bandwidth it is possible to approximate the real-time VBR application to a CBR application with a bit-rate equal to the peak rate of the shaped rt-VBR source. This result suggests that there is an advantage of shaping rt-VBR characteristics toward CBR.

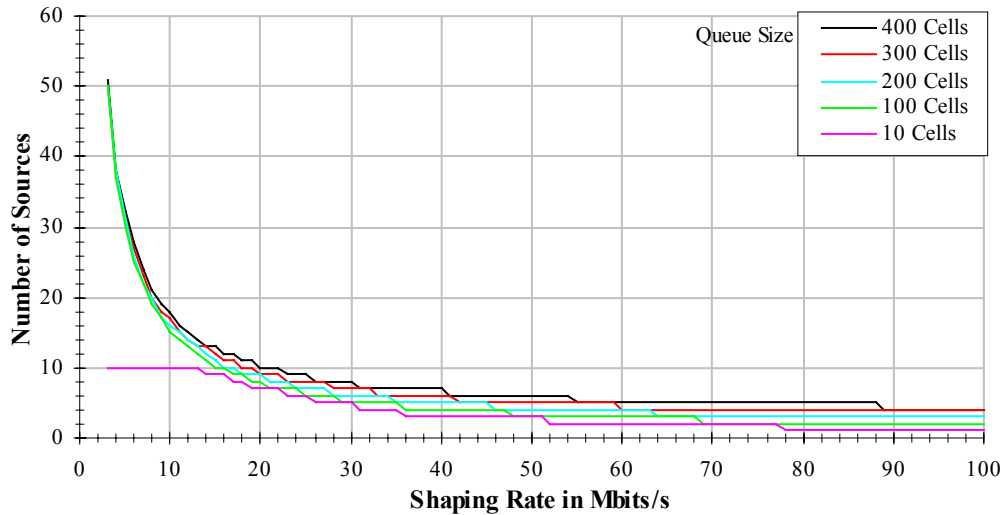


Figure 5.27: ALB for Larger Switch Buffer Sizes using the ISABEL MM Terminal.

The ISABEL model was applied to buffers with larger queues. The largest queue on which simulations were carried out had a maximum queuing delay of 1 ms (400 cells).

Each source had a mean rate of 1.716 Mbit/s. Due to the high peak rate and burstiness of the source, only four sources could be multiplexed when the application's shaper rate was 100 Mbit/s and the queue size was 400. This meant that the network utilisation was extremely low ($\rho=0.044$). Table 5.8 shows the network utilisation for a homogeneous shaped collection of ISABEL terminals. A utilisation (ρ) of 1 means that 100% of the cells on the network are assigned.

Shaping Rate	3 Mbit/s	5Mbit/s	10Mbit/s	15Mbit/s	20Mbit/s	25Mbit/s	50Mbit/s	75Mbit/s	100Mbit/s
Utilisation	0.563	0.353	0.199	0.143	0.110	0.099	0.066	0.055	0.044

Table 5.8: Network Utilisation using ISABEL multimedia terminals

Using a buffer size of 400 cells, the values in table 5.8 demonstrate that even when 100% of the resources are allocated according to the peak rate, there were a large number of empty cells on the network. 43.7% of the network remains unused because the mean bit rate is 1.716 Mbit/s and it is necessary to allocate a PCR of 3Mbit/s. However, this bandwidth will not be lost as non-real-time traffic, ABR (with MCR=0) and UBR traffic can utilise this residual bandwidth to increase the network utilisation to even higher levels.

5.5.4 Summary

Experiment 5 has shown that applying shaping to a traffic stream results in multiplexing density gains without losing end-system QoS. Applying shaping to the end-systems reduces the peak rate of the source; this in turn reduces the bandwidth required and lessens the cost of the rt-VBR circuit, as explained in Experiment 4. Having unshaped sources on the ATM network can cause very low network utilisation. Reducing the bandwidth allocation and burstiness of the end-systems by applying shaping means more sources can be multiplexed on to the network. When the users operate the applications near the QoS “cliff edge” and hence the minimum required CBR bandwidth, the utilisation rises to 56.3% of the network capacity, for a suite of ISABEL terminals.

The user will still have difficulty in characterising the source, even with PCR allocation, as the traffic characteristics depend on many variables. The screen size, audio quality, picture refresh rate, fast moving video scenes and the amount of talk-spurts all cause the traffic characteristic to change substantially. Therefore, the maximum utilisation that was obtained through experimentation is very unlikely when a network of real users attempt to predict the traffic characteristics. However, using the proposed solution of removing rt-VBR category will have gains in simplicity over effective bandwidth calculations and other computationally heavy VBR CAC policies. In addition, by increasing the utilisation of the user’s CBR channel the overall network utilisation increases, as a greater number of users are allowed on to the network, hence increasing the multiplexing density.

6. Discussion

Quality of Service has been defined throughout this thesis as a subjective parameter. This means that the degradation in QoS as seen at the application layer has been affected by every layer on the OSI model and cannot be attributed to a single layer. Therefore, the operator of an ATM network should only guarantee network performance and *not* Quality of Service.

ATM networks have been designed to compete with existing circuit switched networks, while maintaining the apparent unlimited capacity of data networks, like FDDI, SMDS and Ethernet. Therefore, ATM has to be rigorously designed to cope with delay sensitive services and flexible enough to allow IP type services. ATM is a fast hardware switching technology and has the ability to be easily configured into separate networks over isolated virtual paths. ATM is currently being used as an underlying transport network shifting cells quickly from IP entry point to exit points, [ALLE95]. This will enable users to evolve from IP-only networks to native ATM ones.

The proposed QoS mechanism offered within the thesis allows simplicity, stringent guarantees, user defined QoS and efficient use of network capacity. The hypothesis removes the rt-VBR category and relies only on the CBR ATC to provide the real-time service. Using monetary pressure and shaping, use of the real-time service category can be made efficient. CBR applications do not need to be traffic shaped, and from the results of analysis, already offer the most efficient use of a real-time ATCs.

From preliminary experiments, the most prominent VBR traffic type found on ATM networks are from IP services. These applications are naturally very bursty, with long burst lengths and silent periods. The traffic properties of these sources are not fixed, but can be shaped into characteristics that can better benefit the customer and network operator. Generally, rt-VBR traffic has been shown to be very inefficient. Large savings can be made by peak-rate shaping, which reduces the resources needed to transmit each application. This shaped traffic will

improve the multiplexing density on the network through better resource sharing amongst end users, and potentially reduce the complexity of traffic management.

6.1 Cell Delay

When designing an ATM network the most stringent services must be planned and dimensioned: these are the real-time services. The real-time services will be overshadowed by the amount of IP services within the network, so that the switch has to handle at least two types of time priority, real-time and non-real-time. The ATM switch selected must have a high complexity to differentiate between all of the different ATCs on the network, [BOLL96], [BRIE98] and [JIAU96].

The research within this thesis has considered ATM for the most stringent services. The ITU-T has recommended that an ATM switch needs to transmit cells in 300 μ s. From experiment 1, a selection of ATM switches was measured to determine their delay performance. A large range of values was obtained within the selection examined. The fixed delay ranged from 47 μ s to 151 μ s and the maxCTD varied from 140 μ s to 767 μ s. The switches that exceeded a 300 μ s real-time boundary had configurable buffers and were designed primarily for data networks. Using the mean of the fixed delay measurements, an estimate was derived of the nominal size of a real-time queue. Based on real ATM switches with a service rate of 155.52 Mbit/s and a fixed delay of 92.695 μ s, a maximum queue of 76 cells was used as a model to maintain the 300 μ s boundary.

It is difficult to assess the CDV of a traffic stream from the dimensions of a switch, as the effects of CDV are caused mainly by other traffic being multiplexed into the same output buffer. However, there is a maximum boundary for CDV and this is the queue length within the switch. Therefore, the CDV will be dependent on the size of the buffer and the CAC policy used within the network. The multiplexing density of CBR application can be obtained through analysis, [PITT96]. M/D/1 queueing analysis shows that a utilisation of 80% can be obtained for a loss ratio of $1 \cdot 10^{-8}$ with a buffer length of 40 cells.

Within local area networks, the cable delay can be calculated directly from the length of the cable. A propagation delay of 5 μ s per km is usual with optical fibre. When the distances become larger and communication between end-systems takes place through a third party, the exact topology and routing information are not available and approximations have to be made. The ITU-T has proposed a method and the calculation was presented in chapter 3. The ATM Forum states that the network operator should measure the delays with test equipment to determine the CTD and CDV degradation across each link. This is far more accurate than the proposed ITU-T approximation. However, network operators have secrecy policies to withhold the information of network topologies and delays so that an operator can protect the location of customers and cabling route strategies. Thus, the ITU-T method is still the preferred method in the public environment. In experiment 1, the CTD measurement and CDV profile of long distance connections was assessed. A connection across Europe and one across the Atlantic was measured and were found to deviate from the expected result calculated according to the ITU-T. The mean CTD from Switzerland to Holland is 9 ms (ITU-T calculation determined 10.67 ms), and from Switzerland to Canada is 55 ms (ITU-T calculation determined 39.725 ms). CTD needs to be minimised at all times and is heavily dependent on the cabling route and the switches along the end-to-end path. Therefore, CTD needs to be designed into the network when the network is constructed. CTD is a parameter that can be specified at connection set-up, however, the possibilities for varying this parameter are limited. It is in the interest of the network operator to reduce the latency of cells within the network: the longer a cell remains in the network the more expensive it becomes to transmit the cell, as the longer paths need more equipment, additional functionality is needed in the end-systems to absorb CDV and the throughput of TCP/IP is significantly reduced with delay.

CDV has been identified by the ITU-T as a complexity dependent parameter; the more switches transversed, the greater the CDV. However, the measurements of CDV were similar: to Holland, seven switches were crossed to give a CDV of 0.12 ms and for the path to Canada, 13 switches were crossed and a CDV of 0.17 ms was measured. These are much smaller than the ITU-T calculation

suggested, (2.44 ms for each connection). The probable reasons for the low CDV values are that all connections across Europe are peak rate allocated, low levels of utilisation were expected, and the number of independent connections was low over the network. This shows that CDV is as much a measure of interfering traffic as it is the network complexity.

There are many other factors a network operator must take into account when considering the performance of ATM networks, factors such as link synchronisation, background traffic, physical path re-routing and hardware failures. These factors will all contribute to additional CTD, CDV and CLR degradation.

6.2 Effects of CDV on CBR Traffic

N-ISDN telephony services are inherently CBR and range from a single telephone channel of 64 kbit/s to a group of channels e.g. an E1 channel (2.048 Mbit/s). Many video streams have a CBR characteristic, although MPEG-1 and MPEG-2 use VBR encoding, many MPEG devices will smooth the traffic into a CBR profile.

CBR services are affected by CTD, although from experiment 1, the European connection (9 ms) and the transatlantic connection (55 ms) did not present a problem to the CBR applications. Delays of 25 ms start to become problematic with interactive conversations, requiring echo cancellers to be installed. It is generally recommended that the end-to-end delay remain as short as possible.

CDV is a greater problem to CBR services than CTD. CDV is introduced into the traffic stream when traffic is multiplexed together on a shared medium. Queueing analysis is possible for some types of traffic and the most popular models are the M/D/1 and $\Sigma N_i D/D/1$. [DEL20] gives a comparison of analytical solutions to real network measurements. One of the greatest influences on the QoS for a CBR connection is the size of the play-out buffer. Using the long distance connections from Experiment 1 as an example, the peak-to-peak CDV on the connections were 0.12 and 0.17 ms. CDV is dependent on the bit rate being transmitted between the

end-systems; this has not been taken into account by the ITU-T recommended calculation. The ITU-T predicted peak-to-peak CDV for both the European and Transatlantic connections to be 2.44 ms.

The worst case CDV, found through simulation, was dependent on the bit rate and for a 64 kbit/s application was 17.65 ms. A 2.048 Mbit/s application had a peak-to-peak CDV of 3.125 ms and a 39 Mbit/s application experienced 0.918 ms of peak-to-peak CDV. From the results, the CDV expected at the play-out buffer was much larger than on real networks. By provisioning for large CDV, worst case network scenarios have been accommodated.

Compensating for CDV adds additional CTD to the cell stream. From experiment 1, CTD from the network can be quite low (9 ms and 55 ms on the connections tested) and in experiment 2, the CTD added by the receiver contributed a significant reduction in the end-to-end CTD performance; at 210 kbit/s a delay of 18.157 ms was experienced. This meant that the delay caused by the end-systems can be as detrimental to CBR services over ATM networks as the network delay itself. The ATM CBR services need to compete technologically with existing synchronous networks and CTD and CDV should be minimised.

6.3 Application Subjective QoS

The perceived CTD has been shown similar with all interactive applications. Too much delay would cause inconvenience to the speaker due to lengthy pauses between talk spurts and more importantly, data communications is affected in a similar way. TCP uses a system of data transfer and acknowledgements and is very dependent on the delay between the end systems. Therefore, TCP throughput can be significantly reduced if the delay is too high.

As CDV increases, CBR applications suffer from losses in the information stream due to the “empty” and “overflow” states of the play-out queue, causing effects such as loss of sound and picture freezing. As the CDV increased further the applications tended to crash and no communication was possible. Thus, to increase resilience of CBR applications to CDV, longer play-out queues are

needed in the input of the end-system. However, from experiment 2, increasing the queue size increases the effective delay across the network.

Using rt-VBR applications, CDV was examined through the introduction of a shaping device. As the peak rate of the shaper was lowered, the amount of delay accumulated by the burst of traffic increased, the last cell of the burst experiencing the largest delay. Through simulation, the largest individual cell delay was found to be 23.614 ms and the application withstood this amount of peak-to-peak CDV without any visible effects. The probable cause for the robustness of this application was that the processing delay was very large and the CDV was absorbed within the video and audio processing mechanism.

The amount of CLR associated with CBR applications proved to be high. Acceptance levels do not only depend directly on the amount of cell loss, but also on the “expectation” level of the service. Cell loss ratios of $1 \cdot 10^{-5}$ to $5 \cdot 10^{-4}$ are acceptable levels of service for the CBR applications tested. With a 64 kbit/s telephony connection, conversations were still easily understandable with loss ratios of $5 \cdot 10^{-3}$ and only became a little more difficult to understand with $8 \cdot 10^{-2}$ losses. However, the end user of an N-ISDN unit may not necessarily be human. If a computer were attached to the N-ISDN adapter, then the computer’s goodput would worsen considerably with the cell loss ratios presented. The network operator has no way of knowing if the communication is from a human-to-human or computer-to-computer and hence must take into account the worst case scenario and provision for computer communication.

Using the rt-VBR ISABEL application, the CLR was associated to the QoS. It was determined that a CLR worse than $1 \cdot 10^{-4}$ would render the application unusable. The application is based on IP, so when a single cell loss occurred, the application could not reconstruct the IP packet and thus the packet was discarded. This meant that the effective error rate on higher layers was much higher.

To conclude, subjective QoS for a set of applications has been mapped using ATM layer network performance measures. Performance bounds have been determined for four real-time applications. There are limits for CTD, CDV and

CLR. Through experiment, it has been found that, for ATM to compete with other synchronous services, these measures must be minimised. The CTD should be small, allowing greater throughput for applications that need feedback. CDV needs to be minimised so that the magnitude of the processing delay, due to the play-out queue, can be kept to a minimum. Likewise, CLR should remain small so that applications having low loss requirements, particularly computer communications, can sustain high throughputs by maintaining low retransmission rates. Finally, from these high performance ATM circuits users' QoS can be met. The network operator must maintain these circuits in order to provision for all foreseeable services. The users have direct feedback from the applications and so any cost saving made by reducing the performance parameters must be done in the users' premises. The operator cannot make any assumption about the end systems and must dimension for the worst-case traffic types at all times.

6.4 User Benefits of Shaping rt-VBR

Advances in applications have led to a requirement for high quality rt-VBR services, primarily for applications with a video content; applications like multimedia terminals and video servers are primarily rt-VBR sources. One particular application, the ISABEL multimedia terminal, [ISAB97], was used to investigate real-time services. This application is based on IP and needs real-time links to minimise the delay between end-users. The traffic characteristics of the ISABEL terminal could be altered substantially with the introduction of a shaping device. While maintaining the same mean (1.716 Mbit/s), the peak-rate and burstiness of the application was altered without affecting the QoS. The peak rate can range from 100 Mbit/s to 2.2 Mbit/s without the application failing or there being a perceived deterioration in QoS.

Rt-VBR is a complex ATC; the user needs to specify the peak rate, sustainable cell rate and the maximum burst tolerance of the source within the traffic contract. This is almost impossible to specify as the resources needed in terms of sustainable cell rate and maximum burst size can only be discovered once a connection has taken place. This is because the mean rate and burst size may depend on each user's system configuration, speech pattern, video content or

terminal performance. Therefore, it is necessary for the user to make a “best guess” of the parameters in the traffic contract using “rules of thumb”. This guess may be imprecise and the user needs other criteria to determine the best traffic contract that would suit his communication needs. Further criteria that will affect a user’s choice are in the form of charging, which can have a great effect on the types of traffic allowed on to the network. As described in experiment 4, bursty traffic does not promote an efficient use of network resources over real time circuits. Therefore, the unshaped traffic characteristic from an IP based multimedia terminal would need a similar amount of resources as a 100 Mbit/s CBR connection. Applying charges to promote the efficient use of bandwidth will encourage users to reduce the peak rate requirement and alter their traffic characteristics.

Much research has been done on modelling VBR sources to obtain accurate traffic characteristics. Experiment 4 has shown that a source characteristic will vary considerably. A multimedia terminal and a file transfer have been studied and from the results it is seen to be possible to change the characteristics of the sources, using a shaping device, without degrading the QoS perceived by the user. Therefore, using “rules of thumb”, a charging mechanism and a shaping device, the customer can achieve the same QoS performance using cheaper lower bandwidth circuits. This allows the user to select an optimum price / quality function that will make the service cheaper.

6.5 Operator Benefits of Shaping rt-VBR.

In Experiment 5, the multiplexing density from a collection of unshaped multimedia applications proved to be rather poor. The ISABEL application’s mean rate was 1.716 Mbit/s without shaping and had a peak-rate equivalent to the link capacity. To maintain QoS, only 2 multimedia applications were allowed on the 149.76 Mbit/s (STM1) network with a buffer size of 100 cells. The maximum utilisation was 2.2% of the networks capacity because of the bursty nature of the sources and the small queues within the switch. From experiment 1, the switch queue cannot be made larger as it would exceed the constraint of a maximum switch delay of 300 μ s. To increase the utilisation, shaping of the individual

connections was incorporated. From experiment 4, the ISABEL terminal was used with standard user settings, and an optimum amount of peak-rate shaping was 4 Mbit/s. Once this optimum was determined, the statistical multiplexing gain was identified. When optimally shaped, it was possible to multiplex 37 multimedia terminals ($\rho = 0.41$), as opposed to two ($\rho = 0.044$) with a switch queue size of 100 cells. When examining the number of connections and multiplexing density of the VBR connections, comparisons were made with CBR analysis. The comparison proved very close, which indicated that very little multiplexing density is possible with rt-VBR connections. The ability of a small queue to absorb any burst scale queueing is limited. This means that the main criterion for acceptance onto the network is the PCR of the application. There is no advantage for the network operator to provide rt-VBR circuits, the additional complexity and poor multiplexing density make the service uneconomical. From this thesis, it is recommended that rt-VBR traffic should be forwarded through the ATM network over CBR channels. The only possible advantage for rt-VBR channels is to have some estimation of the actual utilisation on the network prior to connection acceptance. This enables an estimation of the residual bandwidth that can be used for lower priority ATCs, e.g. nrt-VBR. However, even without this estimation, nrt-VBR connections can still be transferred over the remaining capacity of the author's proposed QoS mechanism.

6.6 Further Work

The work within this thesis can progress in the future. Analysis of traffic passing through wide area ATM network needs to be completed with large volumes of real time user traffic. This needs to be done to examine a more accurate estimate of CDV on real time circuits. Within this thesis low CDV values were reported on the real networks tested, however, analysis and switch buffer capacities suggest that CDV could increase. Thus, the analytical "worst cases" presented earlier must be adopted until more information on network performance can be gathered. Using the proposed mechanism of traffic shaping rt-VBR circuits, the user can reduce the shaping rate down to the so-called "cliff-edge". Further introduction of background traffic increases the CDV on the network, which causes the shaped

traffic to move over the “cliff edge”, no longer providing acceptable QoS to the user. The present solution to this, is for the user to reduce the CDV in the shaper by increasing the shaper’s service rate, so compensating for the increased network CDV. The proposed further work could examine the sensitivity of application to the increased network CDV while their traffic is being shaped.

7. Conclusion

The research in this thesis has been orientated towards an investigation of the trade-offs involved between QoS and traffic shaping. There are currently two real-time transfer capabilities; CBR and rt-VBR, which have been designed to compete with existing synchronous networks. Due to the nature of ATM networks, CDV is always added to the cell stream. This corrupts the end-to-end timing information needed to synchronise the end-systems' service rate. High performance circuits are needed to maintain this timing relation, so that true synchronous services can be offered over the CBR service category. Therefore, ATM has to be rigorously designed to enable delay sensitive services and flexible enough to incorporate IP type services

There are no native ATM applications available on the market and certainly there will not be any over the short and medium-term time scales. IP is the most important protocol for data communications, and has a large user base that is still increasing rapidly. This type of traffic has been significant throughout this thesis, as the VBR traffic models assume IP over ATM.

The research in this thesis covered the performance of ATM networks, the feasibility of high quality CBR services, QoS to network performance mapping, the user benefit of shaping while maintaining QoS, and the increased network utilisation by incorporating end-system shaping.

The concept of QoS is that of a subjective quality assessment that is dependent on the user's perception of the bearer service. The network operator cannot access the end-systems and, therefore, cannot guarantee QoS. QoS remains the responsibility of the users and the network operator must guarantee the resources necessary for QoS, in terms of network performance.

A flexible QoS framework is needed within an ATM network; it must be simple to use, beneficial to customers and operators, and flexible enough to facilitate both synchronous and IP services. This thesis has proposed a method that provides end users with the required QoS for their applications. The proposed QoS mechanism

allows simplicity, stringent guarantees, user defined QoS, and efficient use of network capacity. The new approach abolishes the rt-VBR category and in so doing, removes the buffering requirements for real-time circuits from the network to the end-systems.

Lower than expected CDV measurements meant that the traffic characteristics did not vary much across the network. This is very important for synchronous services where CDV could cause problems of clock recovery at the receiving side of the network. Low CTD and CDV means that ATM networks, even over long distances, can effectively support stringent real-time services.

A timing relation needs to be maintained across an ATM network so that synchronous services can be offered over the CBR service category. CDV is problematic for this service category, but by suitable provisioning for large CDV worst case scenarios have been accommodated. Compensating for CDV adds additional CTD to the cell stream. This means that the delay caused by the end-systems can be *as* detrimental to CBR services over ATM networks as the network delay itself. The ATM CBR service needs to compete technologically with existing synchronous networks and CTD and CDV should be minimised at all times.

CBR traffic can be placed over an ATM network and coexist with traffic of the same characteristics. In the QoS/Performance trade-off bounded criteria must be met. Applications need minimum levels of performance to give suitable QoS. Research was completed on a set of applications to map the minimum QoS on to network performance parameters. This research showed that CBR ATCs, using a peak-rate CAC policy, could be used to deliver end-to-end QoS.

During experiments conducted on the testbed, the characteristics of several *identical* rt-VBR applications were shown to be quite different. Therefore, predicting the traffic characteristics are difficult and the specification of the traffic contract before transmission is, in most cases, impossible. VBR characteristics are not fixed, but can be shaped into profiles that can better benefit the customer and network operator. Shaping in end terminals means that the buffering capacity for

the source is moved from the network to the user's end-system. The research determined that an application could reduce the PCR parameter without losing QoS. This allows the customers to obtain a cheaper service for the same QoS.

The QoS mechanisms proposed provide a suitable service to users by incorporating: differentiation between ATCs; charging; and shaping. The network operator can apply monetary charges to place economic pressure on users to change their traffic characteristic to facilitate a greater network utilisation. When examining the amount of circuits and multiplexing density of the VBR connections, comparisons were made with CBR analysis. The comparison proved very close, which indicated that very little multiplexing density is possible with rt-VBR connections. From this thesis, it is recommended that rt-VBR traffic should be forwarded through the ATM network over CBR channels. The proposed QoS mechanism simplifies the traffic contract and CAC mechanism, enables even the most stringent synchronous services, and makes efficient use of a prime commodity.

8. Author's Publications

- [BAUM98] M. Baumann, T. Muller, W. Ooghe, A. Santos, S. Winstanley M. Zeller. Multi-Layer Modelling of a Multimedia Application. Broadband Communications '98, Stuttgart, April 1-3, 1998.
- [BRIN96] A.Brinkmann, C. Larvrijsen, S.Louis, R. Macfadyen, D. De Vleeschauwer, R. de Vries, S.B Winstanley. Performance Evaluation - Test Configuration and Results, Telecommunication Systems Vol 5, p249-272 1996.
- [GRIF96] J.M. Griffiths, S.B.Winstanley. Markov Chain Animation Extensions to more than one dimension and time varying inputs, IFIP 4th Workshop on performance Modelling and Evaluation of ATM Networks, Ilkley, July 1996.
- [LOUI94] S.Louis D.D,Vleeshauwer, R.D.Vries,H.V.D. Berg F.V.D. Eijnden S.B.Winstanley, Cell Delay Variation in ATM Networks: Causes and Magnitudes, RACE2064 Traffic Workshop, 14-15 September 1994 Basel.
- [PANK97] F. Panken, J.M. Barcelo, B. Miah, S.B. Winstanley. Investigation on Delay and CDV in an ATM-based Passive Optical Network. IEEE ATM Workshop, Proceedings, 1997, pp. 467-476
- [WINS94] S.B.Winstanley, L.G. Cuthbert Permit Delivery Rates of an ATM Cell Based Access Network MAC Protocol, 11th IEE Teletraffic Symposium, Performance Engineering in Telecommunications Networks Cambridge, 23-25 March 1994, No. 078, pp. 8A/1-8A/9
- [WINS97] S.B.Winstanley, VBR Optimisation in ATM Networks, James User Forum, Invited Paper, Munich September 4-5 1997.

9. References

- [AARS97] E. Aarstad, S. Blaabjerg, F. Cerdan, S. Peters, Kathleen Spaey. Experimental Investigation of CAC and Effective Bandwidth for Video and Data, ACTS094 ATM Traffic Symposium, September 17-18 1997
- [ALLE95] A.Alles, ATM Interworking, Engineering InterOp, Las Vegas, March 1995.
<http://www.cisco.com/warp/public/614/12.html>
- [ASX295] FORE systems, ForeRunner ASX-200BX/ASX-200BXE ATM Switch User's manual, MANU0026-Rev, Warrendale, May 1995.
- [ATM010] ATM Forum Recommendation, ATM User-Network Interface Specification V3.1, AF-UNI-0010.002, 1994
- [ATM032] ATM Forum Recommendation, Circuit Emulation Interoperability Specification, AF-SAA-032.00, September 1995.
- [ATM055] ATM Forum Recommendation, Private Network-Network Interface Specification 1.0, AF-PNNI-0055.000, March 1996.
- [ATM056] ATM Forum Recommendation, Traffic Management Specification Version 4.0. AT-TM-0056.000, April 1996.
- [ATM810] ATM Forum Recommendation, ATM Forum Performance Testing Specification, Doc. ATM96-0810R1
- [BANE97] S.Banerjee D. Tipper, B.H.Martin-Wiess, A. Khalil Traffic Experiments on the vBNS Wide Area ATM Network, IEEE Communications Magazine, August 1997, Vol.35, No.8, pp 126-133
- [BAU96a] M. Baumann, Manual of the simulator: YATS-Yet Another Tiny Simulator, Technical University of Dresden, 1996. <http://www.tu-dresden.de/TK/yats/yats/html>
- [BAU96b] Baumann M.: Results of Traffic Measurements with the ISABEL Application: Construction and Evaluation of a source Model. (AC094_TUD_4.1.1_001.01_CD_CC), Technical University of Dresden, 1996.

- [BELL96] T.E.Bell, J.A.Adam, S.J.Lowe, Communications, IEE Spectrum Journal, January 1996, pp30-41.
- [BLON91] C. Blondia and O. Casals, Cell Loss Probabilities in a Statistical Multiplexer in ATM Network, Proc. 6. GI/ITG Fachtagung Messung, Modellierung und Bewertung von Rechensystemen, Neubiberg, Sept. 1991.
- [BOLL96] R.Bolla, F.Davoli, M.Marcgese, Evaluation of a cell loss rate computation method in ATM multiplexers with multiple bursty sources and different traffic classes, IEEE Global Telecommunications Conference, 1996, Vol.1, pp. 437-441
- [BRIE98] U.Briem, E.Wallmeier, C.Beck, F.Matthiesen. Traffic Management for an ATM Switch with Per-VC Queueing: Concept and Implementation, IEEE Communications Magazine January 1998
- [CHIA96] Leonardo Chiariglione - Convenor, Short MPEG-1 description, ISO/IEC JTC1/SC29/WG11 N1385, June 1996, (http://drogo.cselt.stet.it/mpeg/mpeg_1.htm).
- [COST96] J.Roberts, U.Mocci, J.Virtamo (ed.): Broadband Network Traffic, Final Report of Action 242, Springer Verlag, 1996.
- [CUTH93] L.G. Cuthbert, J.C. Sapanel, ATM, The Broadband Telecommunications Solution, IEE Telecommunications Series 29, May 1993.
- [DANT94] O.M.Danthine P.E.Boyer, Benefits of a Spacer/Controller in an ATM WAN: Preliminary Traffic Measurements, RACE2064 Traffic Workshop, Basel 14-15 September 1994
- [DASI97] L.A.DaSilva J.B. Evans, D. Niehaus, V.S. Frost, R. Jonkman, B.O. Lee, G.Y. Lazarou, ATM WAN Performance Tools, Experiments, and Results. IEEE Communications Magazine, August 1997, Vol.35, No.8, pp 118-124.
- [DEL06] S.B.Winstanley (Editor), Deliverable 06: Specification of Inetgrated Traffic Control Architecture. WP4.1/WP4.2 AC094, September 1996

- [DEL10] S.Peeters, K. Spaey (Editors) Deliverable 10: First Results from Trials of Optimised Traffic Control Features. WP4.1/WP4.2 AC094, March 1997
- [DEL15] S.B.Winstanley (Editor) Definition of Optimum Traffic Control Parameters & Results of Trials, WP4.1 AC094 December 1997.
- [DEL20] S.B.Winstanley (Editor), WP3.3, Deliverable 20, Results of Experiments on Performance Assessment with Artificial Sources, RACE2064 June 1994.
- [DEL30] S.B.Winstanley (Editor), WP3.3, Deliverable 30, Results of Experiments on Performance Assessment using Real Applications, RACE2064 December 1994.
- [DEL31] U.Assmus (Editor), WP1.2, Deliverable 31, Terminal Adapter for High Quality Low Bit Rate Audio Signals, RACE2064, January 1995
- [DEL48] S.B.Winstanley (Editor), WP3.6, Deliverable 48, Final Report on Network Performance, RACE2064 March 1994.
- [DOMI94] J.Domingo , H. Michiel, R. Haberman, M. Sommer, N. Mitrou, Y. Du, M. Gotz, R. Lehnert, Z. Sun, J.P. Cosmas, T. Renger, T. Theimer, Switching Block Studies, Network Performance Evaluation and Traffic Engineering for ATM, European Transactions on Telecommunications, vol.5, No.2 Mar-Apr. 1994 pp.187-198
- [ELWA97] A.Elwalid D.Mitra, Traffic Shaping at a Network Node: Theory, Optimum Design, Admission Control, IEEE INFOCOM Proceedings, 1997, Vol.2 pp 444-454
- [FARB97] J. Farber, M. Frank, J. Charzinski, The WWW-Service in the AMUSE Field Trials: Usage Evaluation and Traffic Modelling. AC094 ATM Traffic Symposium September 17-18 1997
- [G.822] ITU-T Recommendation I.822, Controlled Slip Rate Objectives on a International Digital Connection, 1980 Geneva.
- [G.823] ITU-T Recommendation I.823, The Control of Jitter and Wander within Digital Networks which are based on the 2048 kbit/s Hierarchy, Helsinki 1993.

- [G.824] ITU-T Recommendation G.824, The Control of the Jitter and Wander within Networks which are based on the 1544 kbit/s Hierarchy, Helsinki 1993
- [GRAF97] M.Graf, VBR Video over ATM: Reducing Network Resource Requirements Through End System Traffic Shaping, IEEE INFOCOM Proceedings, 1997, Vol.1 pp 48-57
- [GRIF96] J.M. Griffiths, J.M. Pitts, "Markov chain animation for ATM bandwidth derivation with tandem switching, Performance modelling and evaluation of ATM networks"; Chapman & Hall, 1996.
- [I.350] ITU-T Recommendation I.350, General Aspects of Quality of Service and Network Performance in Digital Networks, Including ISDNs.
- [I.356] ITU-T Recommendation I.356, B-ISDN ATM Layer Cell Transfer Performance, May 1996, Geneva.
- [I.362] ITU-T Recommendation I.362, B-ISDN ATM Adaptation Layer - Functional Description.
- [I.363] ITU-T Recommendation I.363, B-ISDN ATM Adaptation Layer - Specification.
- [I.371] ITU-T recommendation I.371, Traffic Control and Congestion Control in B-ISDN. Geneva, May 1996
- [I.610] ITU-T recommendation I.610, B-ISDN Operation and Maintenance Principles and Functions, Geneva 1995.
- [ISAB97] Home page of the ISABEL CSCW Application <http://selva.dit.upm.es/~proy/isabel/> , December 1997.
- [JIAU96] J.C.Jaiu, C.S.Wu, K.J.Chen, Integrated strategy of resource engineering and control for ATM networks, IEEE Global Telecommunications Conference, 1996, Vol.3, pp . 1739-1743
- [JUNG96] J.I.Jung, Translation of User's QoS Requirements into ATM Performance parameters in B-ISDN, Computer Networks and ISDN Systems, October 1996, Vol.28, No.13, pp 1753-1767
- [KAWA96] T.Kawasaki M. Nakashima, T. Soumiya, M. Katoh, S. Uriu S. Kakuma, A Strategy of Quality Control on ATM

- Switching Network - Quality Control Path (QCP), IEE Global Telecommunications Conference, 1996, Vol.1, pp 432-436
- [KELL96] F.P.Kelly: Notes on Effective Bandwidth, Stochastic Networks: Theory and Applications, Oxford University Press, 1996, pp. 141-168
- [LAND94] R.Landry I.Stavrakakis, A Queueing Study of Peak-Rate Enforcement for Jitter Reduction in ATM Networks, IEEE Global Telecommunications Conference, 1994, Vol.1, pp 619-623
- [LAND95] R.Landry I.Stavrakakis, Traffic Shaping of a tagged Stream in an ATM Network: Approximate end-to-end Analysis, IEEE INFOCOM Proceedings, 1995, Vol.1, pp 162-169
- [LEE96] H.Lee, C.G.Park, Y.H.Kim, Providing multiclass QoS in shared ATM output multiplexer, IEEE Global Telecommunications Conference, 1996, Vol.1, pp 424-431
- [MERA96] Luis A. Merayo, Javier Alonso, and Jesús Marino: Service Clock Recovering in CBR Services: Adaptive vs. SRTS, Proc. IFIP / IEEE Broadband Communications '96, 1996, Montreal, pp. 608-616.
- [MUST96] D.Mustill, Quality of Service Guarantees: A Challenge for ATM Networks, IEEE Global Telecommunications Conference, 1996, Vol.1, pp 420-423.
- [NASE96] H.Naser A.LeonGarcia, Simulation Study of Delay Variation in ATM - Networks, Part I: CBR Traffic, IEEE INFOCOM Proceedings, 1996, Vol.1, pp 393-400
- [NIES93] G.Niestegge, Impact of Traffic Shaping on Broadband Network Performance, AEU, Archiv fuer Elektronik und Uebertragungstechnik: Electronics and Communication, Sep-Nov 1993, Vol.47, No.5-6, pp 426-434
- [ONVU95] Raif O. Onvural, Asynchronous Transfer Mode Networks: Performance Issues, 2nd edition, Artech House, 1995.
- [PITT96] Pitts J.M. Schormans J.A., Introduction to ATM Design and Performance, John Wiley & Sons, 1996
- [RAFF96] Raffali-Schreinemachers, S. Rao, F.Kelly, Y.Markopoulos, Charging and Billing Issues in High Speed Heterogeneous Networking Environments.

- [SBA94] FORE Systems, SBA-200 ATM SBus Adapter Users Manual, Revision Level D, software version 2.3, October 1994.
- [SCHM97] R.M.Schmid, S.Giordano, R. Beeler, H. Flinck, J.Y. LeBloudec., IP over ATM- a position paper. AC094 ATM Traffic Symposium September 17-18 1997
- [SRIK96] A.Srikitja, M.A. Stover, T. Zhong, S. Banerjee, D. Tipper, M.B. Weiss, Analysis of Traffic Measurements on a Wide Area ATM Network, IEEE Global Telecommunications Conference, 1996, Vol.1, pp 778-782
- [VLEE95] D. De Vleeschauwer, S. Wittevrongel, H. Bruneel A. Brinkmann. "Queueing Delay in ATM Switches: Comparison of Measurements and Simulations".
- [WANG96] S.Y.Wang S.E. Meyer, M.Y. Lai, Quality of Service (QoS) for Residential Broadband Video Service, IEE Global Telecommunications Conference, 1996 Vol.3, pp 1949-1953.