

AN INTEGRATED FRAMEWORK FOR ENHANCING USER
EXPERIENCE FOR DIFFERENT DATA SERVICES
IN MULTI-RATE WIRELESS SYSTEMS

by
SOURAV PAL

Presented to the Faculty of the Graduate School of
The University of Texas at Arlington in Partial Fulfillment
of the Requirements
for the Degree of

DOCTOR OF PHILOSOPHY

THE UNIVERSITY OF TEXAS AT ARLINGTON

August 2007

Copyright © by SOURAV PAL 2007
All Rights Reserved

ACKNOWLEDGEMENTS

I would like to thank my Ph.D. advisor Dr. Sajal K. Das for providing the opportunity to work on my doctoral degree at the Center for Research in Wireless Mobility and Networking (CReWMaN). He has been extremely helpful in developing the direction of my research. I have been also fortunate enough to work with Dr. Mainak Chatterjee.

One of the best things that happened to me in graduate school was working with Professor Kalyan Basu. He has been a motivating factor to me in all aspects of life. I am also thankful to my committee members for being kind enough to serve in my doctoral committee.

During the August of 2005, I started working with Sumantra Kundu who is a graduate student at CReWMaN. The collaboration has been extremely fruitful and rewarding. I would also like to thank the other members of CReWMaN who has made my life enjoyable as a Ph.D. student. I would also like to thank my beloved friend Kuver Sinha for being there for me during all the crisis I have faced.

I am also indebted to the Computer Science and Engineering department at the University of Texas, Arlington for funding me during my entire period of graduate study without which I would definitely not have been able to complete my doctoral studies.

Last, but definitely not the least, I would like to thank my parents and my family for providing constant emotional and moral support without which I would not have reached this far. My work, life and existence stand on their sacrifices.

July 16, 2007

ABSTRACT

AN INTEGRATED FRAMEWORK FOR ENHANCING USER EXPERIENCE FOR DIFFERENT DATA SERVICES IN MULTI-RATE WIRELESS SYSTEMS

Publication No. _____

SOURAV PAL, Ph.D.

The University of Texas at Arlington, 2007

Supervising Professors: Sajal K. Das and Kalyan Basu

Technological advances in multi-rate wireless systems have made wireless data services an intrinsic part of human life. An abundance of wireless devices, both in the wireless enabled home multimedia systems as well as in corporate offices, have triggered an array of research in enhancing the wireless data services. The success and acceptance of various rich data services depend on satisfying the user experience derived from such services.

In this dissertation, we first focus on identifying and exposing the important network parameters which hinder user satisfaction in a multi-rate wireless multimedia system. Subsequently, we propose an integrated framework that encapsulates channel state estimation techniques, optimal data rate allocation and user scheduling algorithms, and mobility management solutions. The aim is to improve the performance of data services in multi-rate wireless systems for improving user experience. To this end, we develop an analytical model that correlates user satisfaction with the performance of various user services such as voice, video, and data. Although traditional network performance metrics

such as *throughput*, *connectivity*, and *delay* constraints are also appropriate for multi-rate wireless multimedia systems, the uncertainty posed by the underlying wireless channel poses a new type of challenge in trying to support multiple data rates over the same physical medium. Consequently, methodologies of abstracting the channel state information (CSI), data rate allocation, and user scheduling are different from non multi-rate systems.

In this dissertation, we develop accurate channel estimation techniques specifically for multi-rate wireless systems that enhance the overall throughput of the system. Scheduling algorithms have been specifically designed take advantage of the proposed channel estimation techniques. We have observed that the joint scheduling and channel estimation technique vastly improves the performance of the multi-rate wireless systems in terms of effective throughput and user satisfaction. We have also proposed rate adaptation techniques and content aware scheduling which enhance the performance of streaming multimedia. In the case of throughput constraints, we have focused on determining a theoretical upper bound on the number of satisfied users in comparison to existing schemes for a single cell in a multi-rate wireless system. In addition, we have also focused on mobility solutions specifically for wireless LANs in order to ensure user satisfaction. To ensure connectivity and honor the strict timing requirements demanded by streaming multimedia applications, we have designed and implemented a client-end handoff framework for Wi-Fi systems using the Madwifi driver. All the proposed mechanisms jointly enhance the performance of the data services for multi-rate wireless systems and consequently maximize the number of satisfied user.

TABLE OF CONTENTS

ACKNOWLEDGEMENTS	iii
ABSTRACT	iv
LIST OF FIGURES	x
LIST OF TABLES	xiv
Chapter	
1. INTRODUCTION	1
1.1 Problem Exposition	4
1.2 Challenges Encountered	6
1.3 Related Work	8
1.4 Contributions of this Dissertation	10
1.5 Organization of the Dissertation	13
2. MODELING USER SATISFACTION FOR HETEROGENEOUS WIRELESS DATA SERVICES	14
2.0.1 Contributions of the QoS Framework	15
2.1 Architecture and Policies	16
2.1.1 Network Architecture	16
2.1.2 Policy Management and SLA Issues	17
2.2 Traffic Classes and QoS Metrics	18
2.2.1 Conversation Class	19
2.2.2 Streaming Class	20
2.2.3 Interactive and Background Classes	20
2.3 Modeling User Irritation Factor (UIF)	21
2.3.1 Short Term User Irritation Factor (SUIF)	23
2.3.2 Long Term User Irritation Factor	26
2.4 Radio Resource Management	27

2.5	Simulation Model and Results	32
2.5.1	FER modeling	32
2.5.2	Traffic Models	33
2.5.3	Results	35
2.6	Summary	38
3.	METHODOLOGY FOR ACCURATE CHANNEL ESTIMATION AND SCHEDULING IN MULTI-RATE WIRELESS SYSTEMS	39
3.1	Contributions	40
3.2	Methodology for Channel Estimation in Multi-Rate Systems	41
3.2.1	Problem Formulation	43
3.3	Information Theoretic Framework	44
3.3.1	Why Entropy Measures?	44
3.3.2	Rényi's Error Entropy	46
3.3.3	Proposed Channel Rate Adaptation	47
3.4	Scheduling for Multi-Rate CDMA	49
3.4.1	Throughput Maximization (TM)	49
3.4.2	Fairness Constrained Throughput Maximization (FCTM)	51
3.4.3	QoS and Fairness Constrained Throughput Maximization (QFCTM)	53
3.5	Rate Adaptation in MAC for Multi-rate Systems	55
3.5.1	Fine Grain Scalable (FGS) MPEG-4	55
3.5.2	Rate Adaptation Technique	58
3.5.3	Transition Probability Calculation	59
3.5.4	Analytical Modeling of Rate Adaptation	62
3.6	A Case Study: MPEG-4 over HDR	64
3.6.1	Content Aware Scheduling (CAS)	64
3.6.2	CAS Goodput	69
3.7	Simulation and Experimental Results	71
3.7.1	Simulation Results	72

3.7.2	Results for Rate Adaptation and CAS	79
3.8	Summary	81
4.	THEORETICAL BOUNDS FOR MAXIMUM NUMBER OF SATISFIED USERS IN MULTI-RATE WIRELESS SYSTEMS	82
4.1	Problem Formulation	82
4.1.1	Contributions	83
4.1.2	Preliminaries on Auction Theory	86
4.1.3	Scheduling in Wireless Networks: A Qualitative Formulation	87
4.1.4	Wireless System Model	88
4.1.5	Optimal Scheduling of Wireless Users	90
4.1.6	LP Formulation	91
4.2	Modeling User Utility	92
4.2.1	Utility Function and the Exposure Problem	93
4.3	Scheduling and Combinatorial Auctions	95
4.3.1	Forward Auction for Multi-Rate Slot Allocation	97
4.4	Multiple Slot Scheduling through Reverse Auctions	98
4.4.1	Deprivation Function	99
4.4.2	Mechanism for Reverse Auction	100
4.4.3	Restricted Phase	101
4.4.4	Unrestricted Phase	102
4.5	Performance Analysis	105
4.6	Simulation Study	108
4.6.1	System and Channel Model	108
4.6.2	Simulation Results	109
4.7	Summary	113
5.	EMANCIPATING THE IEEE 802.11 NETWORK FROM HANDOFF DELAY	115
5.1	Motivations for Fast Handoff	116
5.1.1	Contributions	118

5.1.2	802.11 Handoff: Formal Definition	119
5.2	Dynamics of the Beacon Inter-arrival Time	121
5.2.1	Empirical Data Collection and Initial Observation	123
5.2.2	Why is the Beacon Process Non-Deterministic?	125
5.2.3	Beacon Inter-arrival Time Sequence: Short Range Dependent . . .	127
5.3	Framework for fast Handoff	130
5.3.1	IEEE 802.11 State Machine	131
5.3.2	Bootup Phase and Creation of Initial Backup AP List	133
5.3.3	Service Discovery and Maintenance of APs	133
5.3.4	Handoff Execution	135
5.4	Performance Evaluation	137
5.4.1	Experimental Testbed	138
5.4.2	Experimental Results	139
5.4.3	Impact of Buffering	143
5.5	Comparison with Existing Approaches	146
5.6	Summary	148
6.	CONCLUSIONS	149
6.0.1	Future Work	151
	REFERENCES	152
	BIOGRAPHICAL STATEMENT	161

LIST OF FIGURES

Figure	Page
1.1 Overview of the Wireless Data Delivery for Different Applications in Heterogeneous Networks	2
1.2 Experimental Setup for Case Study	3
1.3 Quality at high RSSI	4
1.4 Quality at low RSSI	5
1.5 Picture Quality at no Loss	6
1.6 Picture Quality at 0.1% Loss	7
2.1 Simplified 3G Network Architecture	17
2.2 QoS Requirement for Different Service Classes	19
2.3 Examples of Utility Functions	22
2.4 Web page traffic scenario	34
2.5 Dropping probability for voice calls	34
2.6 Block probability for voice Calls	36
2.7 Average rate jitter for streaming class	36
2.8 Average bounded delay for elastic traffic	37
2.9 Average negotiated delay for elastic traffic	37
3.1 Channel Estimation and error minimization	45
3.2 Original W matrix	49
3.3 Order Matrix	51
3.4 Adpatation Layer for Content Aware Rate Adaptation	58
3.5 Proability Caculation for case I and II	61
3.6 MAC Scheduler and CAS Scheduler	65
3.7 Comparison of CAS Goodput and QFCTM Throughput	71
3.8 Estimation and adaptation for proposed and current Scheme	73

3.9	Magnified view (Time epoch 400 to 600)	73
3.10	Error in the Proposed and PF Scheme	74
3.11	System Throughput for different scheduling algorithms	74
3.12	Mean throughput per user for different algorithms	75
3.13	Standard deviation (SD) of throughput for different algorithms	75
3.14	Variation on the Percentage of Satisfied Users for different algorithms with respect to number of slots in the schedule cycle	76
3.15	Variation of the Percentage of Satisfied Users for QFCTM for different data rates	76
3.16	Fairness for TM algorithm	77
3.17	Fairness for FCTM algorithm	77
3.18	Fairness for QFCTM algorithm	78
3.19	Variation of PSNR with BER using proposed CAS algorithm	79
3.20	The variation of CAS adaptation due to mobility	79
3.21	Improved throughput of I-frames due to CAS.	80
3.22	Error resiliency of CAS and throughput of I-frames.	81
4.1	Illustration of the Schedule Cycle for Multi-Rate Wireless System	89
4.2	Utility curve illustrating the marginal utility of the user as a function of the data received	94
4.3	Throughput vs. Number of Users in the System. Notice how the throughput decreases with increase in D_{min}	109
4.4	Performance of each scheduling scheme measured using satisfied users as a percentage of the total users	110
4.5	Throughput vs. the number of users in the system for different scheduling algorithms	111
4.6	Slot Distribution of Users in Schedule Cycle	112
4.7	System Throughput and Number of Satisfied users vs. α	112
5.1	High level architecture of the proposed framework	122
5.2	Beacon Inter-arrival Time calculated using data collected from our exper- imental testbed (IEEE 802.11a). As the number of clients accessing the wireless medium increases, it is observed that the beacon inter-arrival time	

becomes less deterministic in nature	122
5.3 Beacon Inter-arrival Time calculated from data collected from UTA wireless network. The data was collected on a weekday from IEEE 802.11g network	123
5.4 Beacon Inter-arrival Time calculated from data collected from a commercial hotspot at Gatwick Airport, UK, April 2005. Network was in IEEE 802.11b mode	124
5.5 Dynamics of the beacon processing at the AP. It consist of five main steps: (1) hardware timer raises the beacon interrupt (2) the interrupt is queued in the global interrupt queue (3) interrupt handler processes the interrupt (4) the beacon frame is created inside the driver (5) the frame beacon is transmitted by the wireless chipset	126
5.6 Dynamics of the beacon processing at the Mobile Node (MN). It consists of four main steps: (1) beacons arrive from the wireless medium and are kept on onchip receive buffer. The hardware raises an interrupt (2) the interrupt is queued in the global interrupt queue; (3) interrupt handler processes the interrupt and is moved to the system RAM (4) the beacon frame enters the IEEE 802.11 state machine and is identified	128
5.7 Autocorrelation struture in beacon interarrival	128
5.8 Variation of mean and variance of the beacon interarrival sequence at the MN	129
5.9 IEEE 802.11 Finite State Machine (FSM) present inside the wireless driver	130
5.10 Communication diagram showing the message flow between our proposed framework and the Wi-Fi network	136
5.11 Illustration of the handoff process. This is also the experimental testbed we have used for evaluating our proposed solution AP ₁ , AP ₂ , AP ₃ and AP ₄ shown are a combination of commercial APs (Linksys WRT54G) and standard Linux PCs running MadWifi driver	138
5.12 Graph of Round Trip Time (RTT) vs received packet for packets generated at 20ms	140
5.13 Graph of Inter-Arrival Time (IAT) vs received packet for packets generated at 20ms	140
5.14 Graph of RTT versus received packet for packets generated at 100ms	141
5.15 Graph of packet IAT vs. received packet for packets generated at 100ms	142
5.16 Packet IAT for VoIP stream as experienced during handoff using proposed solution	143
5.17 Packet IAT time for Video Stream as experienced during handoff using proposed solution	144

5.18	Original Image quality	144
5.19	Image quality with Handoff	145
5.20	Effect of Buffering on Packet Loss for different arrival rates	145

LIST OF TABLES

Table		Page
2.1	Statistics for HTTP Traffic	35
3.1	Specification of the files used in our simulation.	78
4.1	Notations Used in auction based scheduling	88
4.2	D_{min} vs. Maximum Number of users	110
5.1	Comparison of Proposed Solution with existing WLAN handoff solutions	146

CHAPTER 1

INTRODUCTION

The broad objective of the dissertation is the study of data services in multi-rate wireless systems. In stark contrast to traditional wireless systems which are capable of detecting only extreme states of the channel (good or bad), multi-rate wireless systems are able to accurately quantize the channel state information (CSI) into different signal-to-noise ratio (SNR) levels. Consequently, depending on the modulation used at the CSI levels, a wide variety of data rates can be supported. Thus, it is possible to support data services like streaming multimedia in multi-rate wireless systems which are able to support data rates as high as 2 Mbps. The aspiration of the dissertation is to (i) understand the underlying mechanism of multi-rate wireless data systems, (ii) expose the drawbacks which hinder the successful functioning of such systems, (iii) identify the specific network parameters which affect user satisfaction and (iv) propose, design and implement solutions which will enhance the working and performance of the multi-rate wireless systems.

The first generation (1G) wireless systems like Advanced Mobile Phone System (AMPS) [84], and Global System for Mobile communications (GSM) [66] could support a maximum data rate of 9.6 Kbps where as in second generation (2G) systems like Code Division Multiple Access (CDMA) [80] the data rate improved to 20 Kbps. Only with the evolution of multi-rate wireless systems from 2.5G onwards, data rates started improving. With third generation (3G) systems like Universal Mobile Telecommunications System (UMTS), 1x Evolution-Data (1xEV-DO) [1], High Data Rate (HDR)/CDMA2000 Systems [2] the maximum data rate is expected to be more than 2 Mbps. In addition, 802.11 based wireless LAN systems though operating in a small range is a multi-rate wireless

system show great promise in supporting higher data rates (as high as 54Mbps to 108 Mbps).

With the evolution of these high speed wireless systems, wireless data services are turning out to be a reality. Real time data services, streaming multimedia, voice over wireless LANs are some of the numerous applications that are gaining popularity. In Figure 1.1, we elucidate the delivery of data services over the high speed multi-rate wireless systems both in cellular as well as 802.11 systems.

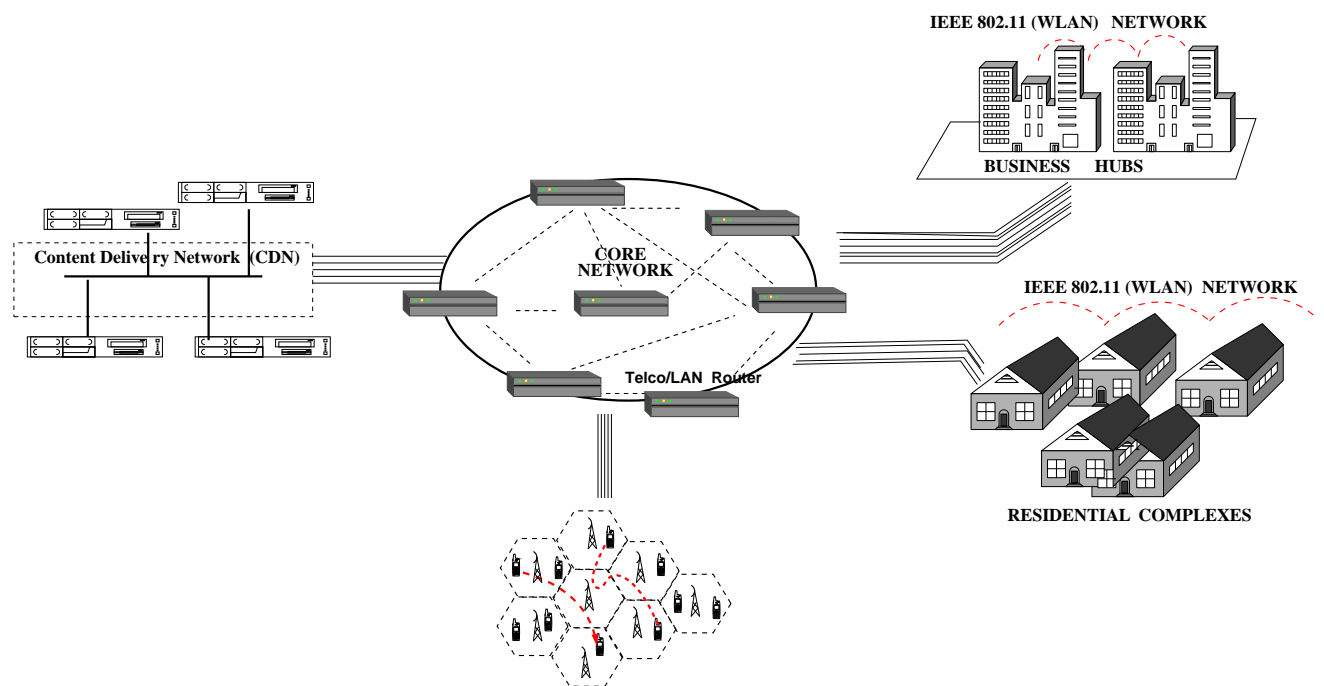


Figure 1.1. Overview of the Wireless Data Delivery for Different Applications in Heterogeneous Networks.

The different categories of applications as illustrated in Figure 1.1 are explained below :

Voice over Wireless LAN (VoWLAN): VoWLAN is a natural extension of VoIP (Voice over Internet Protocol), a technology that has already taken root in enterprise telecommunications. Yet VoWLAN presents its own unique quality of service (QoS)

challenges relating to fluctuating wireless throughput and roaming among APs (access points). The most important QoS demanded is ensuring session connectivity during handoff between two WLAN access point.

Streaming Multimedia over Wireless : Streaming multimedia comprises of both audio and video streams which can be instantly played without the frustration of having to wait for the entire data to be downloaded. Mobile devices which are 3G enabled or having a 802.11 wireless card, are expected to support streaming. However, streaming multimedia has strict timing requirements (maximum possible end-to-end delay is 150ms) and hence suffers from throughput, real time scheduling and connectivity issues when it comes to wireless systems.

The above two applications have emerged as different services in the enterprise market, administrative and academic campuses as well as home entertainment systems. All these services encounter the common challenges of seamless connectivity, throughput enhancement and guarantees for real time scheduling.

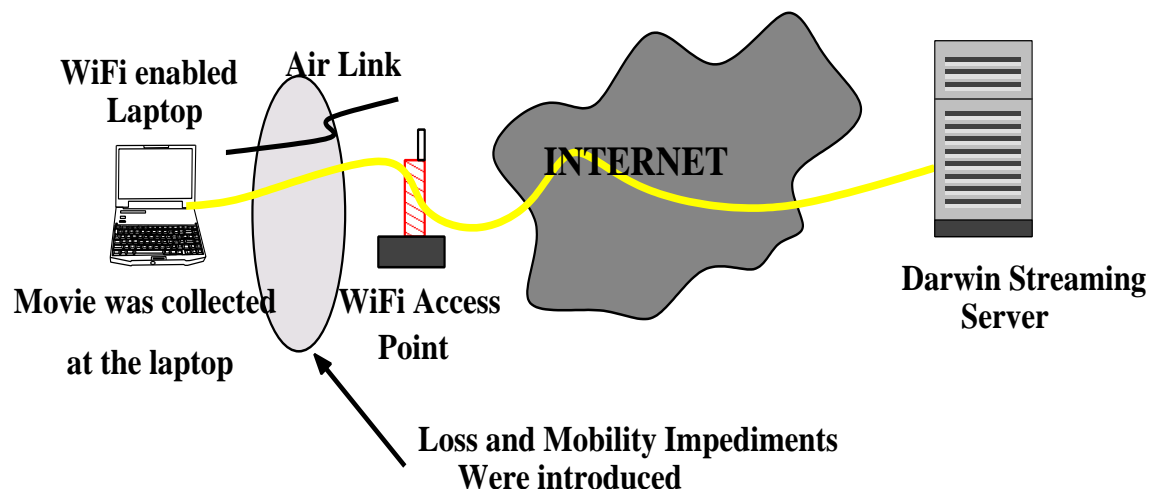


Figure 1.2. Experimental Setup for Case Study.

1.1 Problem Exposition

In order to bring into light the pitfalls of the performance of real time data services of multi-rate wireless systems, we conducted some simple experiments. Before proceeding further, let us examine some of the requirements for the success of streaming multimedia that demand strict timing and synchronization. They are also very much sensitive to network connectivity, delay and loss parameters.



Figure 1.3. Quality at high RSSI.

A small case study exposing how application performance deteriorates due to connectivity as the mobile devices move away from the currently associated access point (AP) justifies the fact that smart mobility management algorithms are needed for maintaining application performance. A Wi-Fi enabled laptop receiving a multimedia stream from a Darwin Streaming Server [23] over 802.11 Access Point is used as the receiving station (STA) for the case study. The experimental setup for the case is shown in Figure 1.2. Initially as shown in Figure 1.3, the video quality is acceptable but as the user walks

to a distance greater than 80 feet the received signal strength indicator (RSSI) falls to a low value of 20 and the video frames start breaking down. It is to be noted that for Atheros chipset AR 5121 [7] the RSSI ranges from 1 to 100 (1 being the worst signal strength indicator). Figures 1.3 and 1.4 elucidate that the STA suffers from degraded data rates and video quality degradation as it moves away from the AP and is in dire need of switching to a new AP which would be capable of supporting the desired required data rate.



Figure 1.4. Quality at low RSSI.

We follow the above with another case study where we introduce loss in the network path of the multi-rate wireless system. The setup used was similar, the only difference was that we specifically introduced a packet loss of 0.01% in the video stream. Instead of playing out the video, we dumped them into corresponding frames. Though it did not result in any loss of frames, it did affect the quality of the movie as can be seen in Figure 1.6. The same frame under the case of network packet loss is shown in Figure 1.5

which conclusively proves that loss still plays a big role when it comes to the quality of streaming multimedia. The multi-rate wireless systems, though capable of streaming media, does suffer from losses and hence are incapable of providing QoS to the media. The same observation was noticed when delay was introduced in the video stream.



Figure 1.5. Picture Quality at no Loss.

1.2 Challenges Encountered

Given the above drawbacks for multi-rate wireless systems, the main challenge is how to model, optimize, and understand the performance of wireless multimedia systems. A lot of work is underway on this broad research area. In this dissertation, we focus on enhancing the user experience for multi-rate wireless systems. Here we enlist the challenges involved :

Modeling User Satisfaction and Resource Management : We envision that the success of wireless data services in conjunction with traditional voice services would



Figure 1.6. Picture Quality at 0.1% Loss.

ultimately depend on *user satisfaction*. Thus a QoS framework needs to be developed that identifies user satisfaction and also facilitates negotiation between the users and the service providers. Identifying the relevant QoS for each of the diverse services and distinguishing the variation of user satisfaction with the perceived QoS is an important research challenge. User satisfaction depends on the subjective expectation of the service and hence varies with the type of services. Modeling the user satisfaction with the varying QoS level for the different services would provide a good insight to designing the resource management schemes for the wireless data systems.

Channel Estimation and Scheduling for Multi-Rate Wireless System : Existing channel estimation and scheduling techniques used for cellular networks do not suffice for the emerging multi-rate systems like CDMA2000, High Data Rate (HDR), WCDMA and wireless LANs. Previous approaches to channel quantization involved just two states - 'good' or 'bad' depending on which channel data rate was decided. However, for multi-rate systems there can be multiple states (8 to 11 depending on the system under

concern) corresponding to varying data rates. Not only does the resource allocation (i.e., scheduling algorithms) need to be modified, the very premise for evaluating their performance need to be changed. We argue that the scheduling performance is very much dependent on the channel conditions and therefore, the scheduling techniques must be strongly coupled with channel estimation.

Connectivity and Fast Handoff in WLANs : Network connectivity in wireless LANs is a big issue when it comes to performance for streaming applications. For any real time applications like VoWLAN fast handoff is a major issue for wireless systems. The inter-packet delay for such applications typically range from 50ms (for example, VoIP) to around 150ms (less demanding audio/video real-time streaming with frame coding at 64Mbps). Applications like Voice over IP (VoIP) or streaming multimedia over 802.11 networks are *real time* wireless data services that *demand seamless* and *continuous* network connectivity. Hence, the challenge is to develop algorithms which guarantee the timing constraints, without breaking the existing 802.11 protocol.

Prior to stating the contribution of this dissertation, we discuss the current state of the art in these fields.

1.3 Related Work

The advent of multi-rate systems like cdma2000 [2], HDR [13], EDGE [29] and 802.11 systems [3] provide higher bandwidth but pose challenges in the joint channel rate estimation and scheduling. The availability of multiple data rates makes the "good-bad" channel modeling inappropriate since accurate estimation of the channel state is necessary to maximize the system throughput. To address these issues in multi-rate systems, several propositions have been made. In [18], the authors have proposed a dynamic rate control algorithm for throughput maximization of multi-rate wireless systems. Yavuz et al. [83] have proposed a hybrid automatic repeat request scheme that improves the

HDR system performance. Very recently in [26], the authors estimated channel condition through a rank-and-subspace tracking algorithm that efficiently adapts the the multipath variations. In [56], the authors propose a channel estimation mechanism based on joint processing of multiburst measurements, relying on the long-term properties of the channel covariance matrix in time-varying propagation environments. The downlink capacity of HDR systems has also been analytically computed in [22] with emphasis on the effect of traffic intensity, radius of the cell, and the scheduling mechanism of assigning slots to active users. Similar results derived from extensive simulations have been presented in [35]. In general, we observe that the throughput of a system is maximized by scheduling the user with the best channel condition, whereas proportionally fair (PF) scheduling guarantees proportional fairness based on the channel state. Both these scheduling mechanisms rely on estimating the data rates for every user. Thus, proper mechanism to estimate the channel state plays a crucial role in deciding the throughput of the system.

Numerous resource management schemes have been proposed which address the various requirements like heterogeneous service demands, fairness, starvation, channel transmission error and delay bounds (for example see [85] and references therein). They include power control algorithms at the physical layer, scheduling or rate-adaptation algorithms at the medium access control (MAC) layer and service admission algorithms at the network layer. However, many such schemes prioritize the voice services and allocate the residual bandwidth to non-real time applications which are deprived of any assurance on the delay bound. A user may care less about the per-packet delay and might put more emphasis on the download time of the *entire* data. This motivates us develop resource management schemes which strives to preserve the QoS requirements of different heterogeneous services, while maintaining fairness amongst different classes of users.

In the literature, significant research has focussed on the scheduling techniques in multi-rate wireless systems addressing varied issues such as user fairness [16], throughput

maximization [13],[35] and efficiency [44]. Existing opportunistic scheduling algorithms exploiting time-varying channel conditions concentrate mainly on throughput maximization while satisfying other QoS requirements. For example, it has been shown in [35], [13] how to maximize the system throughput in CDMA-based HDR systems while maintaining “proportional fairness” amongst users. Similarly, in [44] it has been shown how to formulate the opportunistic problem for multi-channel scenario with resource constraints along with a scheduling scheme that aim to provide fairness among users. The work reported in [40] considered techniques that exploit the wireless channel conditions while guaranteeing each user a predetermined time share in a schedule cycle. In [50], a bandwidth pricing mechanism was proposed that solves congestion related problems in wireless networks. Based on the second-price auction, this scheme shows how the allocation of resources maximizes social welfare. This work was subsequently extended in [49] for designing a pricing mechanism for the downlink transmission power in a CDMA based wireless system. In [75], an auction-based algorithm was proposed that allowed users to compete for time slots in a fading wireless channel. Using the second-price auction mechanism, the users in the system were allocated channel slots and the existence of a Nash equilibrium for such a strategy was proved. Later in [76], the Nash equilibrium strategy was found when the channels for two users are uniformly distributed. To summarize, existing opportunistic scheduling algorithms aim at maximizing the overall system throughput and do not focus on the delay-sensitive requirements of the applications.

1.4 Contributions of this Dissertation

To begin with, we model the user satisfaction - long-term and short-term satisfaction and reveal the main bottlenecks which actually bring down the user experience. In addition, we investigate a QoS framework which respects the user satisfaction constraints. Our findings on variation of user satisfaction with different data services motivate us to

specifically focus on improving the effective network throughput and intelligent scheduling that maximize the number of satisfied users. First we employ an entropy based mechanism to predict wireless channel state and then propose scheduling algorithms for the same. Next, we show that content aware scheduling can greatly enhance the performance for streaming multimedia. We have also used auction based scheduling algorithms as an exercise to figure the theoretical bound of maximum number of satisfied user and how much improvement is possible for the current scheduling algorithms. In addition, mobility introduces connectivity loss which greatly affects user satisfaction. Hence, we have focussed on developing algorithms to ensure fast handoff in 802.11 systems and ensure user satisfaction. In the following, we summarize the major contribution of this thesis :

Modeling User Satisfaction for Heterogeneous Wireless Data Services :

The demand for wireless data services has led to the evolution of third generation (3G) wireless services which deliver a broad range of multimedia applications to mobile users. The first step is to determine how user satisfaction varies with different services. The challenge lies in modeling the subjective user satisfaction to a quantitative metric. In addition, the transition from traditional voice services to data services with heterogeneous requirements necessitates a revisit of the radio resource management schemes. Resource management must consider the impact of error prone transmission medium, heterogeneity of application requirements, and issues related to fairness among users. Also, there is a need to differentiate users based on the amount of revenue they are willing to pay and their expectations from the services.

Methodology for Accurate Channel Estimation and Scheduling in Multi-Rate Wireless Systems :

We used information theoretic techniques to estimate the varying channel conditions as experienced by every user in the system. We utilized a non-parametric estimator

of Rényi's entropy [67] using the Parzen windowing technique to estimate the probability density function of the channel rate variation. Scheduling algorithms are proposed based on the estimated channel conditions that maximize system throughput while ensuring fairness amongst users and maintaining the stipulated data rate requirements. Through a case study (MPEG over HDR), we have shown how packets from a stream are scheduled once the time-slots for a user are determined. Simulations have been conducted and the results demonstrate that the proposed estimation technique coupled with the scheduling algorithm perform better than existing schemes.

Theoretical Bounds for Maximum Number of Satisfied Users in Multi-Rate Wireless Systems :

We investigate the delay-sensitive scheduling problem by borrowing techniques from auction theory [37] and strive to find the maximum possible number of satisfied users that can be served by one access point in a multi-rate wireless system. First we have shown that existing schemes like opportunistic scheduling schemes suffer from a phenomenon that is termed as the *exposure problem*. Then we propose a scheduling technique based on auction theory and demonstrate by mathematical analysis that the proposed scheme is capable of supporting maximum possible satisfied users with hard real time requirements than the existing family of opportunistic and proportionally fair schedulers.. Our approach also leads to significant gain in the system throughput. This study is more of a theoretical finding so as to determine the scope of improvement for existing schedulers in multi-rate wireless systems.

Emancipating the IEEE 802.11 Network from Handoff Delay :

In case of cellular 3G systems, handoff schemes are very robust and developed due to the existence of control channels. Unfortunately in the case of 802.11 systems no such control channels exist and thereby making the handoff algorithms extremely complex. While the IEEE standard defines the specifications for 802.11 Medium Access Control

(MAC) protocol and RF-oriented PHY parameters, it does *not* define any specific roaming algorithm and is open to the device vendors to *improvise*.

Network connectivity is *lost* or the session *severely degraded* if the mobile node fails to associate with a new access point within the time constraints as it moves out from the vicinity of the currently associated AP. We have successfully designed, implemented and evaluated a handoff framework for Wi-Fi networks (IEEE 802.11a/b/g) which achieves a best case delay of 15ms and guarantees that the delay is bounded between 20ms in the average case and 26ms in the worst case.

1.5 Organization of the Dissertation

The dissertation is organized as follows. In Chapter 2, we present our findings on user satisfaction modeling for both long-term and short-term irritation. We also investigate radio resource management schemes which strive to maximise the number of satisfied user. Entropy based channel estimation techniques as well as scheduling algorithms which jointly maximizes the throughput and ensures QoS of each user for multi-rate wireless systems are described in Chapter 3. A novel content based scheduling technique specifically for MPEG-4 over HDR is studied to justify that the proposed techniques result in enhanced performance for streaming multimedia. Thereafter in Chapter 4, we compute the theoretical bounds on the maximum number of satisfied users in multi-rate wireless systems. This is followed in Chapter 5 by the client-end integrated framework which guarantees fast handoff in IEEE 802.11 wireless LANS. Finally in Chapter 6, we summarize the contributions of the dissertation and discuss future work.

CHAPTER 2

MODELING USER SATISFACTION FOR HETEROGENEOUS WIRELESS DATA SERVICES

The success of future generation wireless data services will depend on the parameterized provisioning of quality of service (QoS) for applications whose demands and nature are highly heterogeneous. Also, user satisfaction will play a key role in the economic viability of wireless service deployments both for cellular and Wi-Fi networks. Our objectives in this dissertation are twofold : (i) to authenticate the network parameters which affect user satisfaction, and (ii) to understand how user satisfaction varies with different parameters. with a goal to provide solutions that enhance user satisfaction for wireless data services.

The satisfaction a user derives from the quality of service received is a subjective thing. We deal with this subjectiveness that is involved in user satisfaction or expectation from service providers by defining what we call *user irritation factors*, using Sigmoid functions. These factors reflect users' sensitivity and tolerance to network parameters. To capture the variation of user irritation factors with varying network parameters, we propose a QoS framework based on the paradigm of *traffic class* and *user satisfaction*. The proposed class-based QoS framework comprises a radio resource management scheme which considers user satisfaction based on the perceived QoS, and caters to heterogeneous applications that have diverse QoS requirements. The proposed resource management scheme has two components: the admission control algorithm catering to the *long term* user satisfaction, and the session-based rate and bandwidth allocation scheme manipulating the *short term* user satisfaction. Soft-reservation schemes are also proposed to cater to the higher paying users who would be provided with relatively higher quality of service. Performance metrics have also been specifically identified for each traffic class. In addition to exposing specific network parameters affecting user satisfaction, we demonstrate

that the proposed framework offers improved QoS without compromising the utilization of the system.

The remainder of the chapter is organized as follows. The network architecture along with how the user service level agreement and the policy management issues are stored in the databases are discussed in Section 2.1. The various traffic classes and the right metric for QoS characterization are presented in Section 2.2. The subjective user satisfaction model and its dependence on different service types and perceived QoS is studied in Section 2.3. The two-level radio resource management scheme is proposed in Section 2.4 while Section 2.5 presents the simulation model and the experimental results. Conclusions are drawn in the last Section.

2.0.1 Contributions of the QoS Framework

We propose a radio resource management framework which tries to adhere to the QoS requirements of the applications by exploiting the *subjectiveness* associated with user satisfaction. We categorize users into different classes based on the revenue paid, and consider that all the classes are endowed with heterogeneous wireless services. We propose the notions of *short term irritation* and *long term irritation*, extend them to multiple traffic classes, and propose service (call) admission algorithms and scheduling policies. These policies are tailored for traffic defined Universal Mobile Telecommunications System (UMTS) [1]. Priorities are also considered among these classes. The proposed two-level resource management scheme tries to improve the delay in delivering the entire non-real-time content and real-time traffic (conversational/streaming) by taking into consideration the short- and long-term effects of user irritation. More specifically, the proposed call admission control algorithm regulates the long-term irritation of the users, whereas the short-term satisfaction of the users is guaranteed by the scheduling

policy. It not only provides a bound on the delay but also manipulates the resources so as to maintain the irritation of each user below a certain threshold.

2.1 Architecture and Policies

We present the network architecture along with the databases which contain the user profiles for the proposed QoS framework. These databases also hold the policies pertaining to the user classes and their respective quality of service expectations from the system.

2.1.1 Network Architecture

The resource management framework on which our proposed scheduling algorithm operates is based on the architecture components proposed in the Internet Engineering Task Force (IETF) Policy Information Base (PIB) [27] for differentiated services. Figure 2.1 shows the architectural overview of a 3G wireless cellular network that consists of three important network elements : (1) the airlink and terminal, (2) the base station transceiver system (BTS) and (3) packet data service node (PDSN). Each cell executes a unique copy of the proposed scheduling algorithm for handling requests generated in that cell. In the 3G system, the PDSN includes the gateway functions that interconnect the Internet domain. The PIB includes components like (1) Service Level Agreement (SLA) database, (2) policy decision point (PDP), and (3) policy execution point (PEP). The air link supports *uplink* and *downlink* channels. The uplink channel transmits the request of the clients to the server, where the scheduler schedules the data to the clients through the downlink channels. Though the uplink bandwidth is smaller, we assume that there exists no uplink channel contention between different clients sending requests to the server. Location dependent channel transmission errors which are bursty in nature

needs to be considered for the scheduler design to achieve accuracy and efficiency of the system.

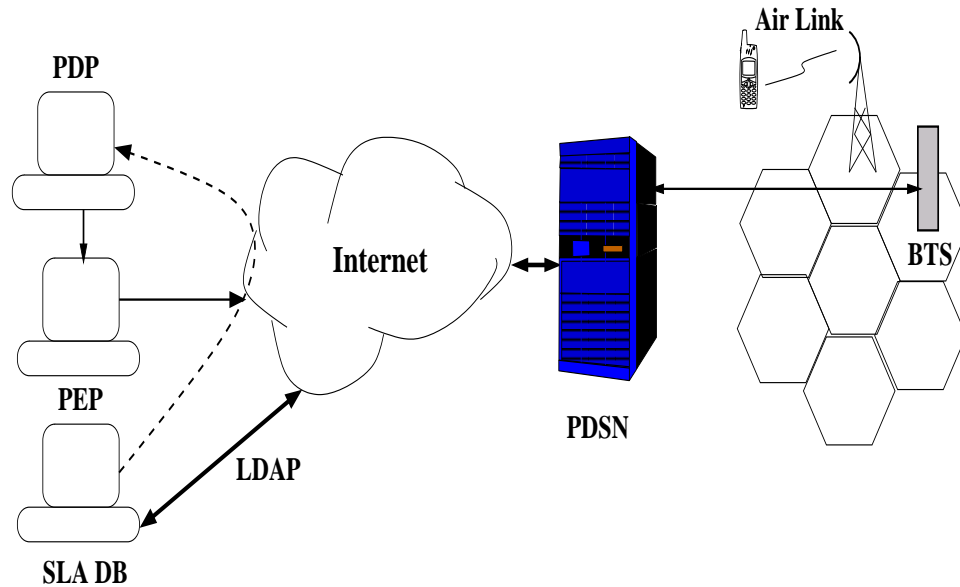


Figure 2.1. Simplified 3G Network Architecture.

2.1.2 Policy Management and SLA Issues

Differentiated service contracts for voice services to mobile users are existant . However, the same cannot be said for data services. The objective here is to create different classes of customers based on the selected service packages for data services. The distinction in QoS levels for these classes lies in the bandwidth provided and hence the delay and throughput offered. The policy management and the customer service agreements are stored in the SLA database as a set of rules. The policy rules are created using the IntServ/DiffServ QPIM (quality policy information modeling) technique as reported in most IETF drafts, for example, in [73]. The PDP contains all the PIBs [27] in addition to the MIBs (management information bases) required for policy management. The PDP function for the differentiated services can be located in the PDSN or mobile switching center(MSC). The policy execution function in this case remains within the radio network

controller (RNC).

We consider three service classes - *Gold*, *Silver*, and *Bronze*, where each user class supports both voice and data services. We propose that the PIB support all three user classes but with different commitment levels. The classification being dependent on the QoS they expect from the network and the revenue they are willing to pay. The Gold class pays the highest revenue and the Bronze class pays the least. The delay suffered for the same service is thus highest for a Bronze class client and least for a Gold Class client. The fairness of service in this context is relative to the price paid by the clients. The scheduler internally classifies the users on the basis of the *user irritation factor (UIF)*, assigning higher UIF to a user with higher priority. This is modeled in the next section.

The PIB also specifies that service to client requests would be processed in the following manner. *Guaranteed QoS* mode of service implies that the system will be able to honor the bounded delay. *Negotiable QoS* is when the system possesses insufficient or no resource at all for a request made. However, if the system is optimistic about being able to serve the request within the bounded delay with the anticipation that on-going transmissions might release resources in near future, then the request is admitted. Nonetheless, the admittance is not strict in nature and the delay restrictions might be violated. The third scenario is when the system is well aware that under no circumstances, it can honor the client request within the specific deadline. It simply rejects the request.

2.2 Traffic Classes and QoS Metrics

The proposed QoS framework caters to the diverse multimedia applications as well as the traditional voice calls in wireless networks. The framework supports multiple classes of users with different priorities and enables fair sharing of the radio resource based on the user class and subjective satisfaction. Moreover, service negotiation between users and service providers as in [61] is flexible in the sense that service classes can be

pre-configured with the user's applications, or explicitly selected at the time of initiation of the applications.

Traffic Class	Conversational	Streaming	Interactive	Background
Fundamental Characteristics	Stringent and Low Delay (<50ms)	Low Delay (< 150ms)	Preserve Payload Contents	Preserve Payload Contents
Streaming	Voice	Streaming Video	Web Browsing	Emails

Figure 2.2. QoS Requirement for Different Service Classes.

For the purpose of illustration and simplicity, we consider the following four QoS classes as proposed for UMTS networks [1, 2]: *conversational*, *streaming*, *interactive*, and *background*. However, the proposed framework is generic enough and can be extended to any number of traffic classes and services as desired by the network operator. These heterogeneous traffic/service classes have specific QoS requirements which has been outlined in Figure 2.2. The main distinguishing factor is the delay sensitivity of each class. Let us outline the different attributes of each traffic class and use them to model the user satisfaction.

2.2.1 Conversation Class

This class is mainly intended to be used to carry real-time traffic flows. Time relation (variation) between information entities of the flow is preserved in conversation class. The other fundamental characteristic is that it has extremely stringent and low

delay requirement. Voice and video telephony are the target applications that falls within the domain of this class.

We use *time-hysteresis outage probability* [86] as the QoS metric for the conversational class. Outage probability is a classical metric in cellular systems: it is defined as the probability that the received signal to noise ratio (SINR) will drop below a specified E_b/N_o , where E_b is the energy per bit and N_o is the noise power. The assumption is that the bit rate requirement and the bit error rates (BER) can be mapped onto an equivalent E_b/N_o .

2.2.2 Streaming Class

The streaming class is very similar to the conversation class but less delay sensitive. The other distinguishing factor is that applications in this class, such as streaming video, are uni-directional (i.e., one-way transport). Though delay and bit error rate affect streaming sessions, *rate jitter* is the most important parameter influencing the quality of such traffic. Hence we employ rate jitter as the QoS metric for modeling user satisfaction, or in other words, for quantifying user irritation for streaming class applications.

2.2.3 Interactive and Background Classes

Interactive and background classes are mainly meant to be used by traditional Internet applications such as WWW, Email, Telnet, FTP, and News. Since the delay requirements of these classes are more slack compared to conversational and streaming classes, they provide better error rate by means of channel coding and retransmission. The main difference between Interactive and background classes is that the interactive class is mainly used by interactive applications like network games and chats, while the background class is meant for background traffic such as emails or web downloading. Responsiveness of the interactive applications is ensured by separating interactive and

background applications. Traffic in the interactive class has higher priority in scheduling than background class traffic, so background applications use transmission resources only when interactive applications do not need them. This is very important in wireless environments where the bandwidth is low compared to wire-line networks.

We observe that the delay a user is willing to tolerate before canceling a session can be considered as a suitable metric for these classes. Though the amount of acceptable delay depends on the particular user, it can be treated as a tunable parameter which varies with the user class itself. An additional constraint is that there should be no data loss.

2.3 Modeling User Irritation Factor (UIF)

The success of our scheduling algorithm lies in modeling the irritation/satisfaction of the user. A basic understanding of the client irritation, i.e., what amount of performance degradation the customer is ready to suffer without complaining, will enable the scheduler to estimate the resources (i.e., number of channels) that have to be allocated to satisfy a particular request. This will directly help in maintaining the delay bound and indirectly help the service provider to control the churn factor [45].

In the following, we propose a method to model the user irritation and present two new metrics: *short term user irritation factor (SUIF)* and *long term user irritation factor (LUIF)*. Each factor signifies different levels of user satisfaction as described below. Qualitatively, *SUIF* measures the delay that the user is ready to suffer before deciding to change or cancel the particular request. *LUIF* determines the tolerance or irritation of the user resulting from continued degradation of service after which the client decides to cancel service completely. For different classes as specified in the SLA, the UIF i.e., the specific *SUIF* and *LUIF* will vary. A high priority user paying higher revenue expecting

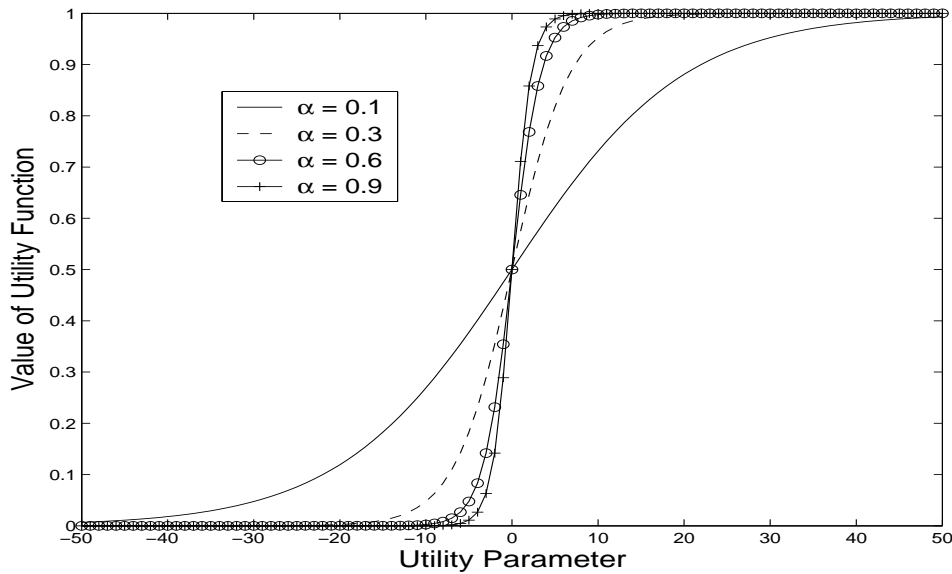


Figure 2.3. Examples of Utility Functions.

lesser delay will be assigned a higher UIF. The goal of the scheduler will be to schedule requests for each client such that the UIF is not violated.

A Sigmoid function has been used in the literature to approximate the user's satisfaction with respect to service qualities or resource allocations [45, 74, 82]. For modeling the satisfaction (or dissatisfaction) of users, we also use the Sigmoid function and correlate it with the proposed metrics, SUIF and LUIF. For a random variable x representing a service parameter like coverage or reliability, the corresponding satisfaction, $U(x)$, is defined by the following sigmoid function:

$$U(x) = \frac{1}{1 + e^{-\nu(x-\beta)}}. \quad (2.1)$$

Here ν and β respectively determine the steepness and the center of the curve and can be tuned to customize the function for different users. The plots for Equation (2.1), for different values of ν are shown in Figure 2.3. From the figure we observe that the satisfaction (utility function) remains zero till it receives a certain specific amount of

service beyond which the satisfaction increases with increasing x . However, after a certain threshold it is observed that the user satisfaction saturates. But for parameters like price or delay, the satisfaction decreases with increasing x . In such cases we can model satisfaction as

$$U(x) = 1 - \frac{1}{1 + e^{-\nu(x-\beta)}}. \quad (2.2)$$

The value of ν indicates user's sensitivity to the QoS degradation while β indicates the "acceptable" region of operation. We use Equation (2.2) to model the UIF. In the following section, we determine SUIF and LUIF analytically for each class of traffic as discussed earlier.

2.3.1 Short Term User Irritation Factor (SUIF)

The SUIF is measured on a per-session, per-user basis. It is also responsible for distinguishing between a call type (both voice and data service) - new or handoff call and accordingly model the user irritation. An in-session user, if deprived of service due to handoff, would suffer from greater irritation than a user whose request is blocked. Hence, we propose a simple mechanism to assign τ_1 and τ_2 signifying the quantitative factors associated with irritation due to a new and handoff call, respectively, where $\tau_1 < \tau_2 < 1$. In the rest of the derivations, we shall use τ to denote τ_1 or τ_2 depending on the request being a handoff or a new call. Also, all the random variables $x_{i,j}$ where i denotes the different service class and j is the user are normalized with the best possible value being 0 (representing zero delay and jitter) and the worst being 1. The value of i will vary between 1 between 4 denoting the four different classes. Equation (2.2) utilizes the $x_{i,j}$ s defined later to measure the SUIFs.

SUIF for Class 1: A classical performance metric in cellular systems is the *outage probability* which is usually defined as the probability that the QoS provided to an existing

connection will drop below a certain threshold, or it is the probability that the received SINR will drop below a specified E_b/N_o . We argue that the SINR translates to the QoS estimation at the physical layer [51]. Classical definitions of outage probability (P_{op}) based on marginal statistics fail to capture the true characterization of the dynamism at the physical layer. It has also been shown in [51] that higher-order statistics of the wireless channel errors affects the performance of the upper layers of the protocol stack. Thus a more general definition of outage probability is needed which considers the *time* dependencies and durations of the unpredictable events. To this end, the concept of *time-hysteresis outage probability* [86] is a more relevant performance metric for voice and data communications. We also feel that it is a good representation and captures the system performance with respect to QoS. Thus, we model the SUIF for class 1 based on time-hysteresis outage probability as the metric for the system performance in our context. If $\mathbf{x}_{1,j}$ denotes the random variable representing the SUIF for class 1 for the j^{th} user, then

$$\mathbf{x}_{1,j} = \tau \times P_{out,j} \quad (2.3)$$

where $P_{out,j}$ is the time-hysteresis outage probability for the j^{th} user. Also $SUIF_{max,V}$ is defined as the threshold SUIF crossing which the voice call is dropped.

SUIF for Class 2: We utilize *rate jitter* as the QoS parameter to model the SUIF for streaming traffic. Jitter, quantified in two ways – *delay jitter* and *rate jitter* is introduced due to variable queuing and propagation delays. Rate jitter [52] measuring the difference between the minimal and maximal inter-arrival packet times, is more appropriate for streaming services than delay jitter. It actually bounds the difference in packet delivery rates during the entire period of service for that particular session and hence is an ideal metric for quantifying user irritation. The higher the rate jitter, the higher is the user

irritation or vice versa. Let $\eta_{max,i,j}$, $\eta_{min,i,j}$, $\psi_{i,j}$ and $\mathbf{x}_{2,j}$ denote the maximum inter-arrival time, minimum inter-arrival, rate jitter and the random variable representing the SUIF for the streaming class for the i th session of the j th user, then

$$\mathbf{x}_{2,j} = \tau \times \frac{\psi_{i,j}}{\eta_{max,i,j}} \quad (2.4)$$

where $\psi_{i,j} = \eta_{max,i} - \eta_{min,i}$ is defined as the rate jitter. Again, $SUIF_{max}$ is defined as the threshold beyond which the call is dropped.

SUIF for Classes 3 and 4: We define the SUIF for the interactive and background classes as the delay that a user is ready to endure before he decides to cancel his request is a measure of his irritation. Additional constraint specific to this class is that there should be absolutely *no* data loss. The ideal transfer time (Δ) for a file of size S Kbytes is $\Delta = \frac{S}{BW}$, where BW Kbps is the ideal bandwidth supported by the system for that user. However, due to congestion, an admitted new request might suffer a delay even before it is scheduled for service. The system can afford to assign a greater delay (δ_{ext}) to class 4 traffic than class 3 since the delay requirements for background traffic is much less strict than interactive. Let the bandwidth demanded by a new request be denoted by BW_d . But the actual delay δ , suffered by the user is given by $\delta = \delta_{ext} + \frac{S_{eff}}{BW_r}$ where BW_r is the bandwidth assigned in reality to the user and S_{eff} is the effective data size [60] that needs to be transmitted due to retransmission on account of frame error rate (FER) instead of the expected $\frac{S}{BW_d}$.

The scheduler is designed to exploit the sensitivity of human nature to delay by transmitting the main page (for Web traffic) and some initial data for class 4 at the earliest possible time. Scheduling the intermediate packets on a regular basis will keep the user satisfied and also provide the scheduler more time to transmit the entire data. We assume that the maximum delay that any user would be ready to tolerate will be

$n \times \Delta$ where n is an integer multiplicative factor. We define the random variable $\mathbf{x}_{(3,4),j}$ denoting the SUIF of the j th user for classes 3 and 4 traffic as

$$\mathbf{x}_{(3,4),j} = \tau \times \frac{\delta_{i,j} - (n - 1)\Delta_{i,j}}{\Delta_{i,j}} \quad (2.5)$$

where $\delta_{i,j}$ and $\Delta_{i,j}$ are the ideal and actual delay for the i th session of the j th user. The worst case bounded delay is $\delta = n\Delta$. and the corresponding SUIF is termed as $\text{SUIF}_{max,BI}$.

For scheduling purpose, we use the *stretch* metric [14] which measures the user tolerance, i.e., additional delay that the user can be made to suffer without crossing the SUIF. We define stretch as the difference between the actual SUIF and SUIF_{max} for the various traffic types.

2.3.2 Long Term User Irritation Factor

The LUIF is a quantitative measure of user tolerance to continuous degradation of the service provided, after which the user decides to cancel the service and churn out of the network. Thus, LUIF keeps track of the long term QoS being provided to each user. The service providers can influence the churn rate by judicious manipulation of the LUIF. Since LUIF is calculated on a per-user basis based on all the SUIFs perceived by that user, maintaining SUIF for each and every request for all users becomes memory and cost intensive. We argue that the QoS received in the distant past would have less significant impact (though not zero) on the users' overall long term irritation than what was received recently. An exponentially weighted moving average (EWMA) mechanism is used to maintain continuous measure of the SUIF for each user. Let the stored LUIF be $U(f_{h-1})$ and the LUIF to be computed be $U(f_h)$; the input to the system, i.e., the current SUIF be $U(x_i)$. The value of $U(x_i)$ can be computed using any one of the

following Equations (2.3), (2.4) and (2.5) depending on the type of request. Then f_h , the random variable which is used to measure the LUIF at the n th request using Equation (2.2) is calculated as follows:

$$f_h = \gamma \times f_{h-1} + (1 - \gamma) \times U(x_i) \quad (2.6)$$

where γ is the weightage given to the cumulative SUIF.

The value of γ needs to be experimentally determined. However, in EWMA mechanisms, $\gamma = 0.2$ or 0.3 is generally chosen. We also define a threshold LUIF signifying that if the LUIF of a particular user exceeds that threshold value, he may cancel the service. This value basically determines the amount of churn in the system. We define another parameter, *tolerance factor*, which measures the amount of patience the particular user has to continue with the current service. Quantitatively, it would be the difference between the threshold LUIF and the actual LUIF. Thus, the call admission control algorithm must always take into consideration that the LUIF of any user does not exceed its threshold.

2.4 Radio Resource Management

We next present the radio resource management scheme offering class-based QoS to heterogeneous services. The proposed scheme consists of two phases: (i) admission control, and (ii) bandwidth reservation and allocation to the admitted requests. The motivation of decoupling the admission control from dynamic bandwidth allocation lies in the premise that admission control is designed to monitor the LUIF of the users, whereas bandwidth reservation concerned with the SUIF determines whether guaranteed or negotiated mode of QoS will be provided. Thus, the QoS framework achieves differentiated service by offering different levels of satisfaction to different classes. Since the framework

strives to adhere to the SUIF and LUIF of different classes, user satisfaction is maximized and hence the churn rate is controlled. The delay bound for a session is computed based on the type of traffic and user class.

The proposed admission control not only admits a higher number of requests by exploiting the delay tolerant nature of elastic traffic, but also endeavors to honor the LUIF of requesting users. The scheme is intelligent enough to judiciously choose and preempt users whose SUIF can be manipulated, or select in-session users who can be delayed/preempted without violating their SUIFs. At the time of a session establishment, the user application specifies its requirements in the form of *Service Request Tuple*. This tuple differs for each class of traffic and is summarized below:

Conversational Class:

$\langle \text{New/Handoff}, \text{BW Required} \rangle$

Streaming Class:

$\langle \text{New/Handoff}, \text{Min/Max BW}, \text{Size} \rangle$

Interactive Class:

$\langle \text{Size} \rangle$

Background Class:

$\langle \text{Size} \rangle$

Although the duration for voice calls is not known apriori for the other request types the ideal delay is computed for making admission decisions.

Algorithm 1 Admission Control Algorithm

```

1: Identify Class, Traffic type of Request
2: Compute SUIF, Retrieve LUIF
3: if Bandwidth Available then
4:   Admit Call
5: else
6:   if Voice/Streaming Request then
7:     Preempt Active Interactive/Background Sessions
8:     if preempted sessions SUIF violated then
9:       Compare LUIFs of requesting and in-session call
10:      if LUIF for requesting call greater then
11:        Drop Call, Update LUIF
12:      else
13:        Admit Call, preempt in-session call
14:      end if
15:    else
16:      Admit Call, preempt in-session call
17:    end if
18:  else
19:    if  $D_r \geq n \times D_i$  then
20:      Admit Call, guaranteed QoS
21:    else
22:      if LUIF for requesting call greater then
23:        Admit Call, negotiated QoS
24:      else
25:        Admit Call, preempt in-session call
26:      end if
27:    end if
28:  end if
29: end if

```

For elastic traffic, in scenarios when the delay bound is violated, the requests are served in the negotiated mode with priority given to higher class users. The admission control algorithm is illustrated in Figure 1. Once a request is admitted, the QoS framework thereafter allocates or reserves bandwidth for that session. We do not allow the entire bandwidth to be accessible to all the classes. This is done to allow Gold class customers to have their deadlines met first. Usually, in case of non-availability of bandwidth, new requests are simply dropped, leading to higher blocking probability. To prevent this situation, a hybrid mechanism involving both preemption and reservation is adopted. We term this as pseudo preemptive service or *soft bandwidth* reservation mechanism. The mechanism is inspired by the following two reasons. Only preemption of lower class clients leads to longer delay for them at the cost of lower blocking probability for higher class clients. Whereas, in case of exclusive reservation of bandwidth for Gold class, the system utilization is low since the reserved bandwidth might be idle at times. Hence the motivation for hybrid approach. We assume a certain fraction of the entire available bandwidth is reserved for Gold class. The reservation scheme can also be extended for Silver class users, but the reserved bandwidth should be less than the Gold class. If the available regular bandwidth is used up by the currently admitted clients, then a lower class client is admitted using the bandwidth reserved for Gold or Silver class. However, the arrival of a higher class request will lead to the preemption of the lower class client from the reserved bandwidth. The detailed bandwidth allocation and reservation algorithm is described in Figure 2. The proposed framework ensures fairness on the basis of the revenue paid since the resource allocation is performed on the basis of user satisfaction as modeled in section 2.3. The lower class users do not suffer from starvation since the admission control scheme takes care of the LUIF of each user. When a request for a user is dropped repetitively, the LUIF of that user is updated by a factor greater than its SUIF.

Algorithm 2 Bandwidth Allocation Algorithm

```
1: if Bandwidth Available then
2:   Admit Call
3: else
4:   if Gold Class then
5:     if Reserved Bandwidth Available then
6:       Serve Call
7:     else
8:       if Lower Class occupies Reserved Bandwidth then
9:         preempt lower class
10:        Serve Call
11:      else
12:        Drop Call, Update LUIF
13:      end if
14:    end if
15:  else
16:    if Reserved Bandwidth Available then
17:      Serve Call
18:      Preempt for Gold class
19:    else
20:      Drop Call, Update LUIF
21:    end if
22:  end if
23: end if
```

This enforces the LUIF to increase towards its threshold value when the scheduler allocates resources so as to prevent the LUIF reaching the threshold. The algorithm takes care of the channel transmission error by considering the effective data size to be transmitted.

2.5 Simulation Model and Results

To validate the two-level resource management framework, we conducted extensive simulation. We evaluate the performance of the proposed class-based QoS framework for heterogeneous wireless services. For simplicity, the simulation considers only a single cell. It was assumed that all 300 users subscribe to heterogeneous services. At any instant, the base station could provide 200 Kbps bandwidth for data services (streaming, interactive, background). The bandwidth reserved for Gold users was set to 10%. Depending on the user class the bandwidth was allocated for data services according to the ratio 4 : 2 : 1. Precisely, Gold, Silver and Bronze class users were assigned 38.4 Kbps, 19.2 Kbps and 9.6 Kbps, respectively. Also, the value of the parameter ν which distinguishes the user class is set to 0.1, 0.3 and 0.9 for Gold, Silver and Bronze users, respectively. For varying system loads, we measured the *blocking probability* for voice calls and streaming requests, *average rate jitter* for streaming requests, *average bounded* and *negotiated delay* for interactive data and background traffic. The effect of frame loss at the link layer and the traffic models under consideration are described next.

2.5.1 FER modeling

Here we describe the modeling of the frame error rate (FER) and its impact on our scheduling algorithm. The dynamically varying capacity of the wireless channel should be sensed by the scheduler and thereafter to take appropriate actions. Although in real life, the packet size varies (if we assume TCP segments), in our design we will assume that the

packet size handled by the scheduler is of constant size. This is made possible using the radio link protocol (RLP) which would fragment a transport layer segment into equal size RLP frames. In order to perform scheduling, the channels are modeled as common pool of bandwidth available for sharing. However, the scheduler must capture the different channel quality (or FER) of each user and the corresponding perceived bandwidth. Due to the nature of wireless medium, packets get lost or damaged, resulting in retransmissions and hence increased bandwidth demand. The effective data that needs to be transmitted is given by $S_{eff} = S/(1 - p)$, where S is the original data size and p is the packet loss rate of the channel. This equality holds true if we do not restrict the number of possible retransmissions. Denoting the FER by p for transmitting n packets, the expected number of retransmissions would be pn . But for the pn retransmitted packets, p^2n more packets are expected to undergo loss or corruption. Thus, due to successive retransmissions, the expected number of packets to be transmitted is given by

$$(1 + p + p^2 + p^3 + \dots)n = \frac{n}{1 - p} \quad (2.7)$$

since $p < 1$. However, for practical systems, the number of retransmissions allowed is finite (usually 3 as reported in [9]). Thus Equation (2.7) can be appropriately modified by considering the required number of finite terms.

2.5.2 Traffic Models

The voice calls are modeled as Poisson process where the inter-arrival time and call holding time are modeled as negative exponential distribution. We modeled the streaming traffic (video) as an *on-off* traffic source, where the *on* time and the *off* time were exponentially distributed with mean value of 30 seconds and 120 seconds, respectively. The traffic for the interactive class was modeled as HTTP [21]. Instead of investigating

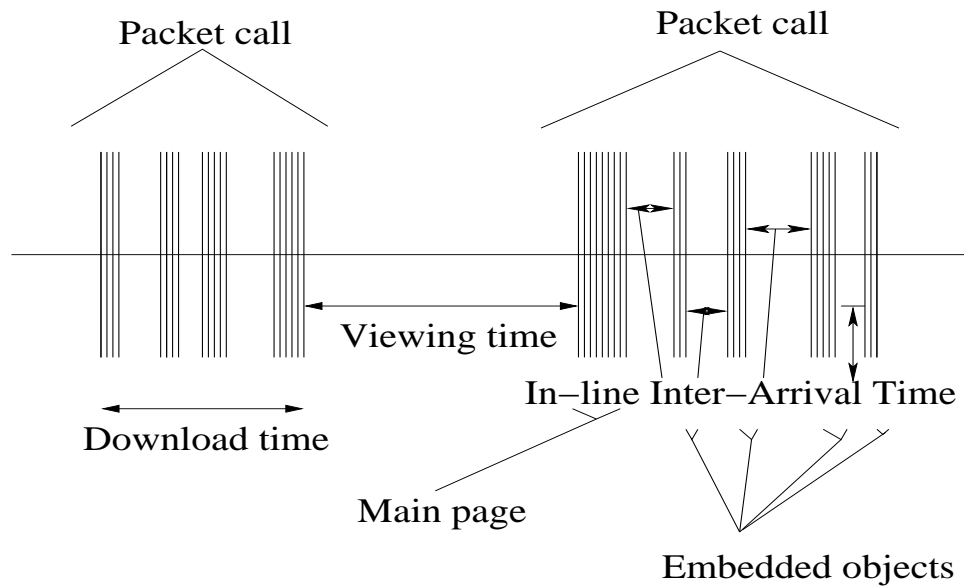


Figure 2.4. Web page traffic scenario.

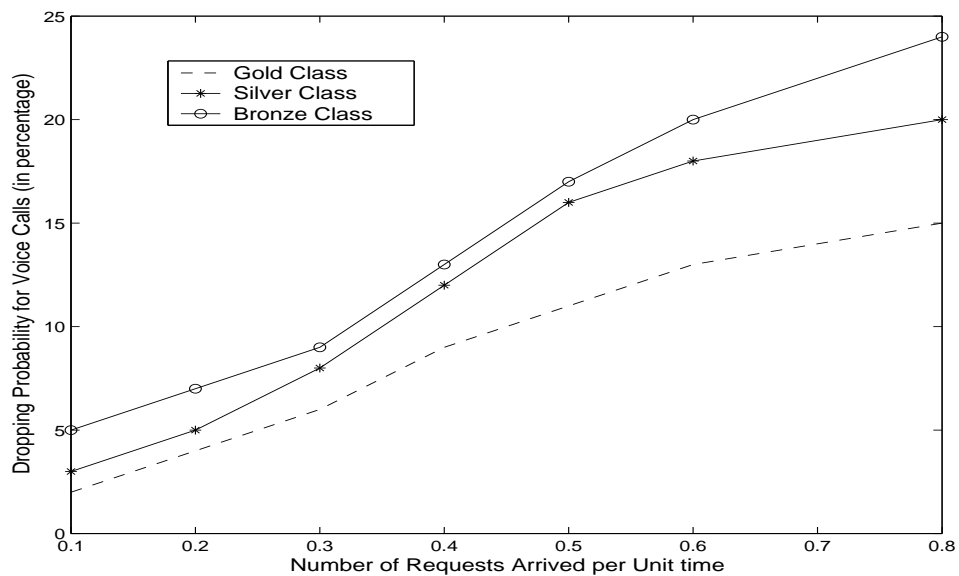


Figure 2.5. Dropping probability for voice calls.

the nature of HTTP traffic, we synthetically generated such traffic with the help of the results obtained in [21]. The basic model of HTTP is shown in Figure 2.4 in which a packet call represents the download of a web page requested by a user. It usually has a main page followed by some embedded objects. A new request (packet call) is

Table 2.1. Statistics for HTTP Traffic

Component	Distribution	Mean
Main page size	Lognormal	10710 bytes
Embedded object size	Lognormal	7758 bytes
No. of embedded objects	Pareto	5.55
Viewing time	Weibull	40 ms

immediately generated after the expiration of the viewing period. The model is also similar to an *on-off* source, where the on state represents the activity of a page request and the off state represents a silent period after all objects in that page are retrieved. The download time of a page follows Weibull distribution, the mean of which depends on the underlying bandwidth of the wireless channel. Each object (main page and embedded objects) of the HTTP traffic is fragmented into multiple equal-sized frames so as to fit into a packet. Other statistics and parameters used to generate the HTTP traffic are shown in Table 3.1. The FTP requests are similar to the web traffic but has only one embedded object and the packet size modeled as a Pareto distribution having different scale and shape parameters.

2.5.3 Results

The dropping and blocking probability for voice calls belonging to all three user classes are shown in Figures 2.5 and 2.6, respectively. Simulation results prove that the QoS received for Gold class is the best followed by Silver and Bronze classes. Thus, the notion of fairness based on revenue paid is adhered to. The rate jitter shown in Figure 2.7 was averaged for the particular class of users so as to measure the average rate jitter for each class. Here also, the Gold class suffers from the least jitter since the bandwidth assigned is maximum for these users. The bounded delay suffered by the background traffic for each of the classes is shown in Figure 2.8. The delay suffered in negotiated QoS

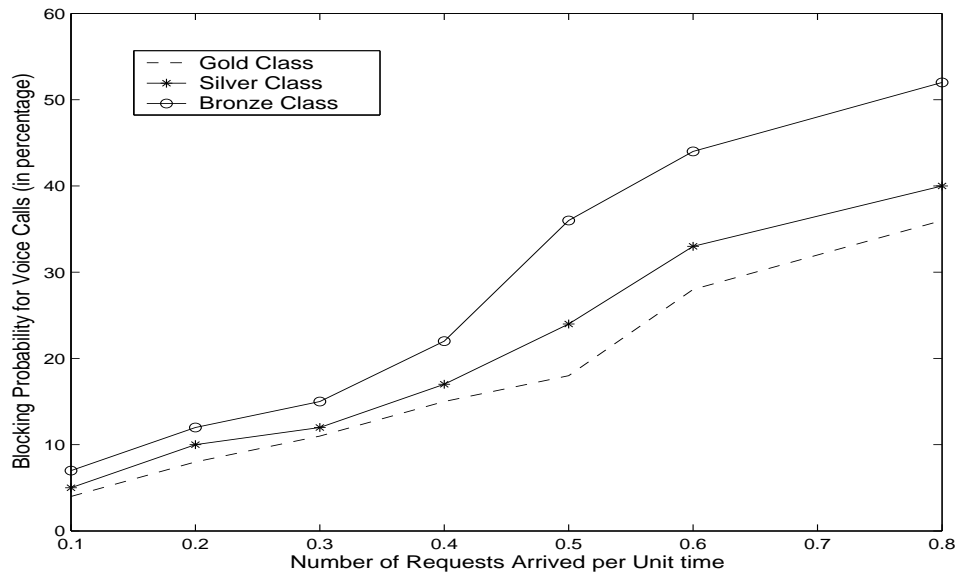


Figure 2.6. Block probability for voice Calls.

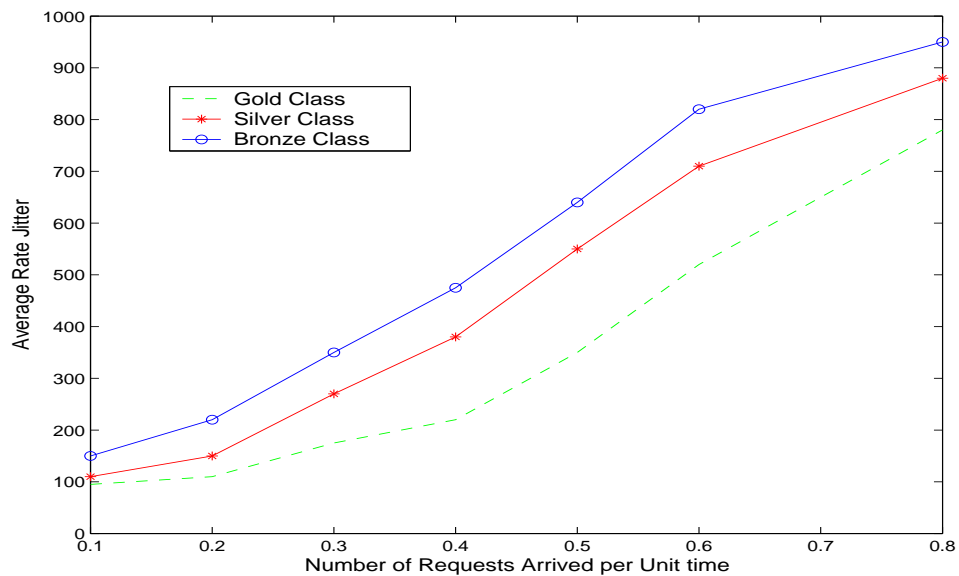


Figure 2.7. Average rate jitter for streaming class.

by the users when the scheduler optimistically admits the user, is presented in Figure 2.9. Note that the delay suffered during negotiation is higher than when served in guaranteed QoS mode. The difference between the delays for the Bronze and higher classes is much higher in the negotiated mode since the background class is penalized more when the

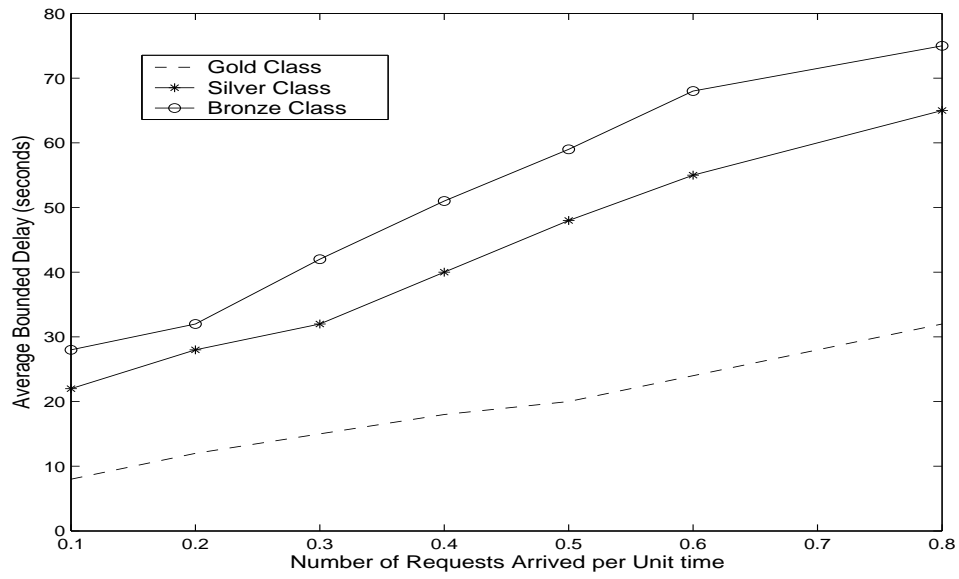


Figure 2.8. Average bounded delay for elastic traffic.

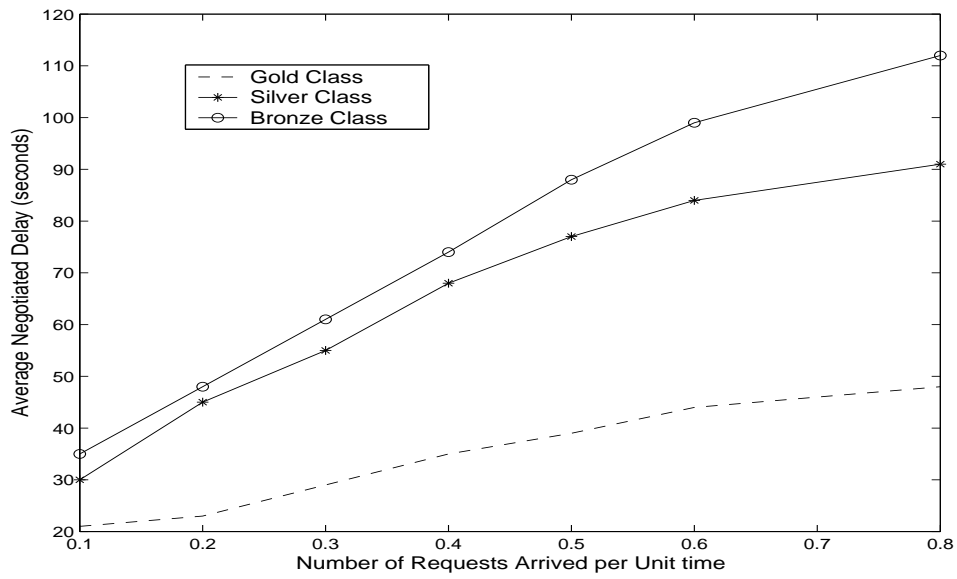


Figure 2.9. Average negotiated delay for elastic traffic.

system is more overloaded. The Bronze class suffers the maximum delay whereas the Gold Class suffers the least, which proves the validity of the proposed class-based QoS framework.

2.6 Summary

A QoS framework for both traditional voice communications and rapidly emerging data services has been proposed to evaluate specifically the network parameters affecting the different services. The framework is based on the traffic class and user satisfaction. Though non-real time wireless data services have elastic requirements, the basic form of QoS or SLA is still non-existent. We provided an insight to what might be the possible QoS parameters for these data services. Variation of user satisfaction to the different services received has been modeled which forms the basis for the QoS framework. The proposed radio resource management schemes comprises of call admission control algorithm and bandwidth allocation policies. The admission control algorithm admits sessions based on not only the resource requirements but also the irritation level of the requesting user. In addition, the framework judiciously penalizes users such that the user irritation factors remain bounded. Simulation results prove that the framework successfully exploits the flexibility in user tolerance with respect to the perceived QoS and provides assurance for non-real time applications. Thus, the results reveal that the proposed algorithms provide bounded delay guarantees and acceptable call blocking and dropping probabilities.

The insights and conclusion we derive are the following : (i) fast and accurate estimation of channel state is needed to enhance the performance wireless data rates and consequently would enhance throughput and user satisfaction (ii) session termination due to handoffs cause higher irritation to user experience; we have focussed on providing seamless mobility solutions for Wi-Fi users (iii) in order to gain insight for performance of base stations in terms of number of active users it can support we have performed a theoretical analysis of the same based on current estimation of channel estimation. We also discuss possible scheduling algorithms which would give better performance in terms of numbers of satisfied users than existing scheduling algorithms.

CHAPTER 3

METHODOLOGY FOR ACCURATE CHANNEL ESTIMATION AND SCHEDULING IN MULTI-RATE WIRELESS SYSTEMS

The findings in the last chapter dictate that new mechanisms are needed to enhance wireless data services in terms of throughput maximization, delay and loss minimization to improve user satisfaction. In this chapter, we focus on the throughput aspect for current multi-rate wireless systems. Our findings reveal that simply maximizing the throughput would not result in maximizing the number of satisfied user. The objective is to support not only the voice calls which require fixed amount of resource but also data services which demand varying resources for maintaining a minimum level of quality of service (QoS). The higher level of delay toleration, bursty traffic characteristics and asymmetrical nature of the data services demand for efficient and judicious *scheduling* of the radio resource. In addition, mobility of users, various propagation losses, and the hostile wireless environment adds to the uncertainty of the channel state. Thus, the first step in providing successful data services in next-generation multi-rate wireless networks comprises of estimating, predicting, and adapting to the channel conditions and thereby ensuring high data rates. The second step comprises of developing scheduling algorithms that would assure per user data requirements thereby ensuring user satisfaction. In addition, fairness and throughput maximization issues need to be considered as well. We have addressed the above two objectives by estimating the channel condition by employing information theoretic techniques and consequently proposed simple but effective scheduling schemes satisfying the above objectives.

The rest of the chapter is organized as follows. Section 3.2 presents the working principle of the existing multi-rate wireless systems and identifies the bottleneck of such existing schemes and formulates the objective to be achieved. A channel rate estimation technique employing information theoretic metrics is proposed in Section 3.3.

Scheduling algorithms guaranteeing general processor sharing (GPS) and assuring data rate fairness criterion based on the estimated channel rate are presented in Section 3.4. A rate adaptation scheme based on the estimated channel throughput is proposed in Section 3.5. Section 3.6 presents a case study of transmission of MPEG-4 over HDR using the proposed techniques. In Section 3.7, comprehensive simulation results involving channel estimation techniques, scheduling algorithms, rate adaptation and performance of MPEG-4 transmission over HDR are presented. Conclusions are drawn in the last section.

3.1 Contributions

In this chapter, we formulate channel estimation problem in multi-rate wireless systems and thereafter maximize the throughput for such systems by employing Rényi's error entropy [25] minimization technique. The probability density function of the error is captured through the Parzen windowing technique [63]. We develop algorithms which guarantee the processor sharing fairness constraint and then further refine it to assure the QoS demands of each user assuming that the demands are known a priori. Although we consider HDR as the reference system, the approach is general enough to be applied to any multi-rate system. In summary, the main contributions are as follows.

- A non-parametric estimator using Parzen windowing technique is employed to measure the online entropy for the multi-rate channel errors. Utilizing Rényi's error entropy minimization technique, the channel state for each slot and each user is estimated.
- We develop scheduling algorithms that maximize the number of satisfied users, guarantee processor sharing fairness, assure data rate fairness, and maintain QoS constraints based on the estimated channel data rates.

- Simulations and numerical analysis have been conducted for the proposed technique. Results for channel estimation techniques and scheduling algorithms which demonstrate significant performance enhancements in terms of number of satisfied users.

3.2 Methodology for Channel Estimation in Multi-Rate Systems

We consider a single cell of a multi-rate wireless system with the base station serving K mobile terminals/users. We assume that the system employs data rate control mechanism on the forward link that adapts to the changing channel conditions by employing adaptive modulation and coding techniques, hybrid Automatic repeat-request (ARQ), and best serving sector selection. The mobile terminals perform measurement of the current channel conditions (i.e., the received $\frac{E_b}{N_o}$) and predicts the achievable rate. Every mobile terminal updates the base station of the predicted rate via the pilot signal on the reverse link data rate control (DRC) channel. At any time slot t , the data rate that can be supported by the i th mobile terminal is $R_i(t)$, $1 \leq i \leq R(N)$, where $R_i(t)$ is one of the $R(N)$ rates supported by the system. For example, HDR supports 11 data rates [13]. $\bar{R}_i(t)$ is the mean rate actually provided to user i measured over a sliding window of length t_c . It is recursively defined as :

$$\bar{R}_i(t+1) = \left(1 - \frac{1}{t_c}\right) \times \bar{R}_i(t) + \frac{1}{t_c} \times R_i(t_c) \quad (3.1)$$

This cumulative estimation is put into effect at each slot; however, the scheduling step is executed once every new transmission.

On receiving the feedback signals on the reverse link, the base station executes the scheduling algorithm to determine the users to be scheduled. The most common scheduling disciplines used in multi-rate systems like HDR are the *maximum Data Rate*

Control (DRC) rule and the *proportionally fair (pf) rule* [13]. The mechanism to decide on which user to schedule by the maximum DRC rule is given by

$$s = \arg \max \{R_i(t)\} \quad (3.2)$$

It implies the user with the highest rate is scheduled. In case of a tie, the scheduler resolves by selecting the user with more data in the queue. The decision for proportionally fair scheduler is given by

$$s = \arg \max \left\{ \frac{R_i(t)}{\overline{R}_i(t)} \right\} \quad (3.3)$$

where $\overline{R}_i(t)$ is computed as per Equation (3.1). Other complex scheduling disciplines like the *modified largest weighted delay first (M-LWDF)* [13] and the *exponential rule* have also been proposed that have sophisticated scheduling decision rules. The M-LWDF scheduler is given by

$$s = \arg \max \{\zeta_i R_i(t) W_i(t)\} \quad (3.4)$$

where $\zeta_i > 0$ and $W_i(t)$ is the wait time of the head of line packet of user i . Similarly, the exponential rule is given by

$$s = \arg \max \zeta_i R_i(t) \exp \left(\frac{a_i W_i(t) - \overline{aW}}{1 + \sqrt{\overline{aW}}} \right) \quad (3.5)$$

where $a_i > 0$ and $\overline{aW} = \frac{1}{N} \sum_i a_i W_i(t)$ However, all these schemes suffer from increasing the system penalty which is demonstrated next.

3.2.1 Problem Formulation

The DRC rule achieves throughput maximization whereas the PF scheduler attains proportional fairness. Moreover, the M-LWDF and the exponential scheduler achieve throughput optimality. However as evident from Equations (3.2),(4.1.3),(3.4), and (3.5), all these schedulers achieve varied objectives related to throughput and stability of the system. The objectives of these schedulers do not relate to the QoS guarantees of the data streams. In comparison, the earliest deadline first (EDF) considers deadline but is not aware of what the achievable rates are for various users in the next few slots, and thus fails to deduce an efficient schedule. We try to overcome the inadequacies of the existing schemes and devise mechanisms to support data streaming exploiting the channel state characteristics, scheduling and adapting the rate of the streams in response to the channel data rate.

Ideally, unit length schedule cycle maximizes throughput. However, such schedulers have no apriori knowledge of the channel state information (CSI). Lacking future CSI, unit length schedulers might fail to guarantee data requirements of all the active users. Preferably, an infinite schedule length with perfect CSI and data requirements will ensure scheduling mechanisms that are fair and meet QoS requirements of all sessions (users). Hence estimation and prediction of the CSI is essential. However, for any prediction technique there is a prediction error associated which increases monotonically with the predicted length. We do not focus on finding the optimal predicted length in this paper but concentrate on devising mechanism to support and maximize the number of data services with QoS guarantees. *The understanding that the channel data rate as well as the data requirements keep varying for each user is fundamental to designing the mechanism to achieve the objectives stated earlier.* We envision that the framework should involve: (i) CSI estimation and prediction and (ii) the scheduling policy should be self-adaptive to the time granularity of the channel variation and data requirement.

In addition, we assert that the proposed approach is ideal for slow fading users, since for fast fading users the channel impulse response changes rapidly within a time slot; hence estimation and prediction techniques fail. In particular, we address the following interlinked questions.

Question 1: How to dynamically estimate and adapt to the channel rate accurately?

Question 2: How would the slots per decision cycle be allocated to mobile terminals such that the QoS requirements are met?

First, we elucidate techniques for online calculation of the channel data rates using information theoretic metrics. Thereafter, we propose three scheduling schemes based on the estimated channel states.

3.3 Information Theoretic Framework

We envision that capturing the *uncertainty* of the channel condition using an information theoretic approach would be the first step to answer the above questions. The ideal measure for uncertainty is *entropy*. We elucidate techniques for real-time calculation of the uncertainty and thereby estimating the channel data rates.

3.3.1 Why Entropy Measures?

The mean square error (MSE) has been a popular criterion for training of adaptive systems [81]. MSE is based on second order statistics i.e. it is able to extract information successfully for linear systems whose statistics can be defined by their mean and variance. For nonlinear systems, entropy as an optimality criterion extends MSE. This is because when entropy is minimized, all moments of the error probability density function (pdf) is minimized. Hence, Shannon's entropy or the more generalized Rényi's entropy [67] are scalar quantities which provide a better measure for the average information contained

in a given probability distribution function. We utilize the principle of minimizing error entropy as the performance criterion in adapting to the channel rate.

Minimizing the Rényi's entropy of the error results in minimization of the divergence between two pdfs [25]. The fact that minimizing the Rényi's error entropy minimizes the Csiszar distance [25] between the pdfs of input-desired signals and the input-output signals has already been established in [25]. The application of the Rényi's error entropy modeling for multi-rate CDMA/HDR system is conceptually straightforward. The variation of the data rates due to the fluctuation of the channel condition is modeled as a non-linear system. The goal is to dynamically identify the non-linear system state error defined as the difference between the desired output and the actual output. For a given set of input samples or channel rates $R_i(t)$ and the output under consideration i.e., the decision rate $R_i^d(t)$, the entropy of the output error $e_i = |R_i^d(t) - R_i(t)|$ must be minimized. Figure 3.1 illustrates the system which captures the error samples and employs Rényi entropy minimization principle to adaptively tune to the channel state of each mobile terminal. Let us now compute the error estimates starting with the definition of Rényi entropy.

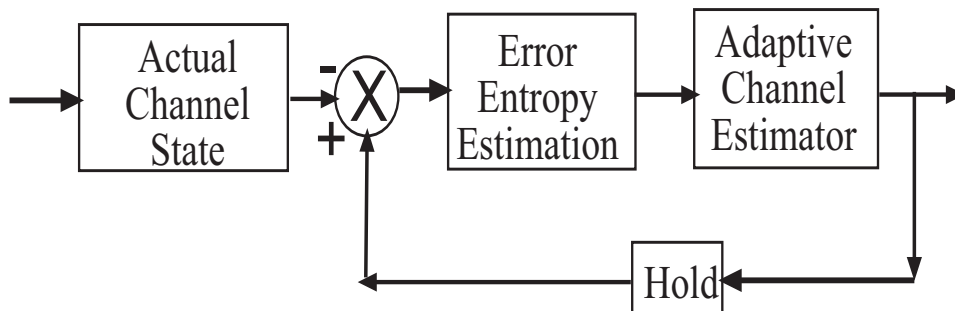


Figure 3.1. Channel Estimation and error minimization.

3.3.2 Rényi's Error Entropy

Rényi's entropy, $H_\alpha^R(P)$, with parameter $\alpha > 0$, for a random variable x with pdf $f_x(\cdot)$ is defined by

$$H_\alpha^R(x) = \frac{1}{1-\alpha} \log \int_{-\infty}^{\infty} f_x^\alpha(x) \quad (3.6)$$

The pdf $f_x(\cdot)$ of the random process in practical cases is unknown a priori. But a non-parametric estimation of the pdf is necessary for evaluating the Rényi's entropy as defined in Equation (3.6). We employ the non-parametric estimator for Rényi's entropy proposed in [24] using Parzen windowing technique [63] with Gaussian kernels to estimate the pdf of the channel rate variation. In Parzen windowing, the pdf is approximated by a sum of kernels whose centers are translated to the sample points. The pdf estimate of a random variable x for which the samples $\{x_1, \dots, x_n\}$ are given, is obtained using the kernel function $\kappa_\sigma(\cdot)$ where parameter σ is the size of the kernel. It is given by

$$f_x(x) = \frac{1}{\Psi} \sum_{i=1}^{\Psi} \kappa_\sigma(x - x_i) \quad (3.7)$$

where Ψ is the number of samples. In the current context, the samples are the channel rates $R_i(t)$ measured in Ψ slots. Thus, using Equations (3.6) and (3.7), we compute the Rényi's error entropy for actual channel rate and the estimated channel rate as follows:

$$H_\alpha^R = \frac{1}{1-\alpha} \log V_\alpha(x) \quad (3.8)$$

where $V_\alpha(x) = \int f_x(x) dx$ is the *information potential*. Since log is a monotonic function, minimization of H_α^R corresponds to maximization of V_α for $\alpha > 1$ or $\alpha < 1$.

3.3.3 Proposed Channel Rate Adaptation

We utilize the information potential as the criterion to adapt to the desired channel rate. The information potential is expanded using the expected value operator as

$$V_\alpha(e) = \int f_e^\alpha(e) de = E[f_e^{\alpha-1}(e)] \approx \frac{1}{\Psi} \sum_i f_e^{\alpha-1}(e) \quad (3.9)$$

Substituting the Parzen window estimator from Equation (3.7) in Equation (3.9) we get the information potential estimator as

$$V_\alpha(x) = \frac{1}{\Psi^\alpha} \sum_j \left(\sum_i \kappa_\rho(x_j - x_i) \right)^{\alpha-1}. \quad (3.10)$$

The kernel function in Parzen windowing is κ_ρ where ρ denotes the width of the window in terms of a predetermined unit-width kernel defined as

$$\kappa_\rho(x) = \frac{\kappa}{\rho} \left(\frac{x}{\rho} \right) \quad (3.11)$$

Thus, we have outlined the methodology for computing the pdf of the channel rates from an incomplete sample set. Utilizing the pdf generated from the Parzen windowing technique, we derive the Rényi's entropy which gives us a quantitative measure of the uncertainty of the channel data rates. The information potential, $V_\alpha(x)$, is maximized based on the error estimates and thereafter adapted to the actual channel rates.

The error samples for a multi-rate system can be computed as $e_i = R_i^d(t) - w_i R_i(t)$ where w is the weight vector that we want to calculate. For initialization purpose, w can be set to 1. The objective is to set w in such a fashion that the error e_i is minimized. On obtaining the error samples, we use the information potential estimator in Equation (3.10) to compute the information content. Thereby, to estimate the optimal

weight vector w^* for a decision cycle of N slots, we obtain the gradient of the information potential estimator as

$$\frac{\partial V_\alpha}{\partial w} = \frac{\alpha - 1}{N^\alpha} \sum_j \left(\sum_i \kappa_\rho(e_j - e_i) \right)^{\alpha-2} \times \left(\sum_i \kappa'_\rho(e_j - e_i)(x_i - x_j)^T \right) \quad (3.12)$$

Using Equation (3.11) and $\kappa'_\rho(x) = \frac{\kappa'}{\rho^2} \left(\frac{x}{\rho} \right)$ we simplify Equation (3.12) to

$$\frac{\partial V_\alpha}{\partial w} = \frac{\alpha - 1}{\rho^\alpha N^\alpha} \sum_j \left(\sum_i \kappa_\rho(\Delta e_w^{ji}) \right)^{\alpha-2} \times \left(\sum_i \kappa'(\Delta e_w^{ji})(x_i - x_j)^T \right) \quad (3.13)$$

where Δe_w^{ji} s are given by

$$\Delta e_w^{ji} = \left(\frac{[(R_j^d(t) - R_i^d(t)) - w^T[(R_j(t) - R_i(t))]]}{\rho} \right) \quad (3.14)$$

From this relationship we can infer that to achieve the optimal w^* , we need to set the weights such that the transpose of w satisfies

$$w^T = \left(\frac{R_j^d(t) - R_i^d(t)}{R_j(t) - R_i(t)} \right) \quad (3.15)$$

Each of the w vectors correspond to one user. The estimate is performed for each of the U users currently admitted in the system and at each decision period, a matrix $W_{U,N}$ is obtained which specifies the data rate for each user at every time slot.

3.4 Scheduling for Multi-Rate CDMA

Given the channel rate estimates it is necessary that the scheduling algorithm chooses the slots according to the desired throughput and QoS demanded while maintaining fairness among users. Since channel rate information for each slot and each user is estimated prior to scheduling decision, multi-user diversity gain can be achieved with the help of proper scheduling algorithm which exploit the channel state information. We present three scheduling algorithms and illustrate with the help of an example as shown in Figure 3.2. The example considers 3 users (rows) and a decision interval of 9 slots (columns) where each slot value signifies the estimated data rate in Kbps. We start with a simple objective of putting into practice the throughput maximization criterion and in steps fine tune the algorithm to accommodate multiple criterion.

user 1	76.8	614.4	614.4	1228.8	921.6	153.6	153.6	76.8	38.4
user 2	1228.8	1843.2	921.6	1228.8	2457.6	307.2	614.4	76.8	76.8
user 3	921.6	38.4	614.4	1228.8	2457.6	614.4	76.8	204.8	2457.6

Figure 3.2. Original W matrix.

3.4.1 Throughput Maximization (TM)

The TM algorithm is similar to the PF algorithm and maximizes the system throughput. However, the PF scheduling imposes some fairness by choosing the user with the highest ratio of the requested to the estimated rate. It is to be noted that the estimated rate is calculated using Equation (3.1). However, the TM scheduling makes use of the estimated data rates for each user for every slot and is presented in Figure 3.

For TM we are only concerned with the throughput maximization and not fairness. We impose fairness and QoS constraints later and revise TM.

Algorithm 3 Throughput Maximization algorithm

- 1: $i = 0, \sigma_{TM} = 0$
 - 2: **while** $i \leq N$ **do**
 - 3: $\sigma_{TM} = \sigma_{TM} + \max(w_{i,j}) \forall j$ for $1 \leq j \leq U$
 - 4: Update SV
 - 5: $i = i + 1$
 - 6: **end while**
-

At every slot, the TM scheduler chooses the user with the highest rate. The throughput for the decision period is given by σ_{TM} . For the W matrix as generated in Figure 3.2, the Schedule Vector (SV_s) is $SV_s = [2, 2, 2, 3, 2, 3, 2, 3, 3]$. The elements in the vector correspond to the user who gets the slots. SV_s is obtained from the TM algorithm where i and j correspond to the slot and the user chosen, respectively. Since the TM algorithm only tries to maximize the throughput, user 1 is starved.

Claim I : *Based on the estimated channel rate the algorithm TM maximizes the system throughput.*

Proof. We prove the above claim by contradiction. Let there exist some other scheduling discipline Γ which provides system throughput σ_Γ higher than TM scheduling (σ_{TM}). Hence $\sigma_\Gamma > \sigma_{TM} \Rightarrow \sum_i^N w_{i,j}^\Gamma > \sum_i^N w_{i,j}^{TM}$ where $w_{i,j}^\Gamma, w_{i,j}^{TM}$ is the set of slots (j) and users (i) chosen by Γ , TM scheduling disciplines respectively. The above relationship implies that there exists at least one $w_{i,j}^\Gamma$ which is greater than the corresponding $w_{i,j}^{TM}$. But this

violates the TM scheduling paradigm which chooses the slot for the user which has the maximum value. Hence the proof. \square

3.4.2 Fairness Constrained Throughput Maximization (FCTM)

The concept of fairness from the scheduler's viewpoint is best captured by the widely known General Processor Sharing (GPS) theory. GPS states that for each user k , there exists a weight ϕ_k , such that the scheduler schedules at least ϕ_k fraction of server time. For a stable system, $\sum_k \phi_k \leq 1$. To incorporate the GPS fairness, we maintain a weight vector $WV[U]$ of length U whose each value corresponds to ϕ_k of the k th user. It is reasonable to assume that ϕ_k 's for all the users are known a priori. If all users have the same priority level then we can conclude that $\phi_k = \frac{1}{U}$. Therefore, each user will have $\lfloor \frac{1}{U} \times N \rfloor$ slot share in a decision period of N slots. The floor function ensures the stability. However, according to priority the ϕ_k 's would be different and WV would be updated accordingly. Let us illustrate the above with 3 clients and consider the decision period limited to 9. We set $\phi_1 = \frac{1}{2}$, $\phi_2 = \frac{1}{3}$, $\phi_3 = \frac{1}{4}$. Hence $WV = [2, 3, 4]$. We order in ascending order each row of W , i.e, slots of each user are ordered for the particular decision period and generate an Order matrix as shown in Figure 3.3 for the W matrix shown in Figure 3.2.

user 1	4	5	2	3	6	7	1	8	9
user 2	5	2	1	4	3	7	6	8	9
user 3	5	9	4	1	3	6	7	8	2

Figure 3.3. Order Matrix.

The Fairness Constrained Throughput Maximization (FCTM) algorithm is presented in Figure 4.

Algorithm 4 Fairness Constrained Throughput Maximization

```

1:  $i = 0, \sigma_{FCTM} = 0,$ 
2: Initialize Order[ $U$ ][ $N$ ],  $WV[U]$ .
3: Reorder in ascending order each row of  $W$ 
4: Update Order[ $U$ ][ $N$ ] with the ordered slot numbers corresponding to the ordered  $W$ .
5: for  $i = 1$  to  $N$  do
6:   for  $j = 1$  to  $U$  do
7:     Choose slots from Order[ $i$ ][ $j$ ]
8:     if ( $!WV[j]$ ) then
9:       if slot Order[ $i$ ][ $j$ ] already not chosen then
10:         $\sigma_{FCTM} = \sigma_{FCTM} + (w_{i,j})$ 
11:        Update  $SV, WV$ 
12:       else
13:        Choose slot with the max value.
14:        Update  $\sigma_{FCTM}, SV, WV$ .
15:       end if
16:     end if
17:   end for
18: end for

```

The Schedule Vector (SV) for the matrix in Figure 3.3 will be $SV_s = [2, 2, 3, 3, 2, 3, 1, 1, 3]$. Thus, FCTM guarantees GPS criterion. However, the throughput (σ) is penalized.

Claim II: *FCTM guarantees GPS fairness criterion and thereafter maximizes system throughput based on the estimated channel rate.*

Proof. The first part of the claim that FCTM guarantees GPS is evident. The number of slot allocation to user is restricted according to the WV vector. We proceed to prove the second part of the claim by contradiction. Let there exist some other scheduling discipline Γ which provides system throughput σ_Γ higher than FCTM scheduling (σ_{FCTM}). Hence $\sigma_\Gamma > \sigma_{FCTM} \Rightarrow \sum_i^N w_{i,j}^\Gamma > \sum_i^N w_{i,j}^{FCTM}$. Note FCTM schedules user j for slot i s.t. user j 's slot rate is maximum for the i th slot and $WV[j]$ is not violated. Thus $\sum_i^N w_{i,j}^\Gamma > \sum_i^N w_{i,j}^{FCTM}$, implies that there exists $w_{i,j}^\Gamma$ which violates $WV[j]$ requirement. Hence proved. \square

3.4.3 QoS and Fairness Constrained Throughput Maximization (QFCTM)

The generalized processor scheduling (GPS) constraint in a wireless domain does not ensure that the data rate demanded by the user is guaranteed. The data rate constraint is much more alluring from the users' perspective. However, it should be pointed out that no *guarantee* can be provided. The data rate demanded by a user can be honored with a probability. This is due to the inherent unreliable wireless transmission medium. In addition, depending on the service (such as voice, multimedia, elastic traffic) being provided, there exist service specific QoS constraints. The ARQ mechanism in HDR allows the base station an update no sooner than 4 slots [35]. We denote this constraint by γ_k for the k th user. The scheduling algorithm should take into account service specific QoS requirements other than the GPS and data rate constraints for each user.

We assume that for each user k , there exists a desired data rate per decision period (DR_k) which needs to be satisfied. We incorporate flexibility in QFCTM scheme such that if DR_k is satisfied for a particular user, then the scheduler overrides the user's ϕ_k requirement. However, the fulfillment of ϕ_k does not necessarily satisfy DR_k . Under such circumstances if there exist extra slots, then the scheduler should schedule those slots to the users whose DR_k is not satisfied. Thus, QFCTM exploits the multi-user diversity in the sense that DR_k for user k might be satisfied without using the full quota of ϕ_k since the user might experience good channel state. The extra slots generated in the above manner are utilized to satisfy the requirements of other users suffering from poor channel conditions. QFCTM maintains a constraint vector $CV[U]$ and a demand vector $D[U]$ besides utilizing the Order and WV as in algorithm FCTM algorithm. Here $CV[i]$ and $D[i]$ correspond to the γ_i and DR_i requirement of the i th user. The QFCTM algorithm is presented in Figure 5.

We illustrate the above algorithm using the same GPS constraints as before and consider data requirements (in bytes for this decision period) of 2000, 1000 and 3000 for users 1, 2 and 3 respectively. We also assume that user 1 cannot be scheduled in successive slots. Then, the desired Schedule Vector (SV_s) for QFCTM is given by $SV_s = [2, 1, 3, 3, 3, 3, 1, 2/3, 1]$. Note that the scheduler is free to grant the 8th slot either to user 2 or user 3 although the DR_i requirement for both of them is satisfied since user 1 cannot be granted any slot due to its QoS requirement.

Claim III: *QFCTM maximizes throughput based on the estimated channel rate and guarantees GPS fairness, and assures data rate and QoS criterion.*

The proof is similar to Claim II.

3.5 Rate Adaptation in MAC for Multi-rate Systems

With the channel rate estimation and scheduling algorithms it is important to devise mechanisms for data streams to be adapted to the changing channel conditions. The link layer aware rate adaptation techniques [10, 12, 20] using statistical measures fail to capture the content level information which is to be exploited for multimedia applications. Moreover, there have been experiments that prove the slow response of TCP with respect to channel rate variation. Thus, a *lag* develops. This lag can be minimized if the adaptation is done at the MAC layer and scheduling is made content-aware. It is known that better rate control algorithms are obtained if stochastic channel behavior through apriori models are considered [6]. Thus, we incorporate the channel model as well as the content information in our rate adaptation algorithm at the MAC layer. The scheduler not only transmits more data than necessary (avoiding buffer overflow) when the conditions are favorable, but also drops low priority frames without compromising on the integrity of the application concerned. The scheduling performed by the MAC scheduler is *content aware* in the sense that the content is classified into multiple priority levels. We would perform a case study MPEG-4 video on our proposed solution to justify the concept of classification based on content. We assume the content being classified as frames. We use content and frame with the above meaning in mind. In order to appreciate the rate adaptation technique presented here a brief background of the MPEG-4 video is needed.

3.5.1 Fine Grain Scalable (FGS) MPEG-4

MPEG-4 [36] is an ISO/IEC standard for MPEG-4 video which is suitable for the harsh wireless environment [72]. It offers flexible encoding bit rate and high compression efficiency by using object-based paradigm of atomic audio visual objects.

Algorithm 5 QoS and Fairness Constrained Throughput Maximization Algorithm

```

1:  $i = 0$ ,  $\sigma_{FCTM} = 0$ ,
2: Initialize Order[ $U$ ][ $N$ ],  $WV[U]$ ,  $CV[U]$ ,  $D[U]$ .
3: Reorder in ascending order each row of  $W$ 
4: Update Order[ $U$ ][ $N$ ] with the ordered slot numbers
5: corresponding to the ordered  $W$ .
6: for  $i = 1$  to  $N$  do
7:   for  $j = 1$  to  $U$  do
8:     if  $D[j]$  satisfied then
9:       Update  $WV[j]$ 
10:    end if
11:    Choose slots from Order[ $i$ ][ $j$ ]
12:    if ( $\neg WV[j]$ ) then
13:      if slot Order[ $i$ ][ $j$ ] already not chosen and  $CV[j]$  satisfied then
14:         $\sigma_{QFCTM} = \sigma_{QFCTM} + (w_{i,j})$ 
15:        Update  $SV$ ,  $WV$ ,  $D$ .
16:      else
17:        Choose slot with the max value.
18:        Update  $\sigma_{QFCTM}$ ,  $SV$ ,  $D$ ,  $WV$ .
19:      end if
20:    end if
21:  end for
22: end for
23: if (slots exist) then
24:   for all  $j$  such that  $D[j]$  not satisfied do
25:     Schedule  $j$  such that  $CV[j]$  is satisfied
26:   end for
27: end if

```

The data units are organized in several hierarchical layers as follows: *Video Session* (VS), *Video Object* (VO), *Video Object Layer* (VOL) and *Video Object Plane* (VOP). The encoded video stream comprises of a Group of VOPs (GOV) structure consisting of I (Intra-coded), B (Bi-directionally coded) and P (Predictive coded)VOPs (frames in case of MPEG-1/2). Thus, low bit rate (as low as 5Kbps), scalable video coding providing Variable Bit Rate (VBR) output, in-built error resilient/concealment techniques make MPEG-4 a very attractive proposition for multimedia wireless transmissions.

However, the basic MPEG-4 standard offers limited granularity in terms of scalability to varying bandwidth requirements with respect to the enhanced fine grained scalable (FGS) video coding. Because of its simple scalable structure, the FGS encoding scheme is capable of covering a *continuous* bandwidth range. The FGS structure has two layers: a base layer (BL) and an enhancement layer (EL). Discrete Cosine Transform (DCT) is used to code both the base and enhancement layer. However, the DCT coding in base layer uses run-level[†] coding whereas the enhanced layer uses *bit-plane coding*[‡]. Any number of bits from the enhancement layer that are transmitted successfully can be decoded and adopted to improve the video quality due to the bit-plane coding. The enhanced stream can thus be truncated at any bit position according to the available bandwidth and thus enables continuous rate adaptation. Detailed description of FGS and its enhancements can be found in [34]. Thus considering the suitability of the FGS scheme for rate adaptation we use FGS encoded video stream for multimedia transmission over wireless.

[†]The number of consecutive zeros before a non-zero DCT coefficient is called a *run* and the absolute value of a non-zero DCT coefficient is called a *level*

[‡]Bit-plane coding represents each DCT coefficient as a string of binary bits.

3.5.2 Rate Adaptation Technique

Successful transmission of I-frames is more crucial than others. This knowledge of frame priority at the MAC layer helps the rate adaptation. The adaptation layer is shown in Figure 3.4 This is because the delay in retransmission of I-frames from the application layer is higher than retransmission from the MAC layer. In addition, the MAC layer can drop less important frames based on the priority and closely follow the channel conditions by minimizing the Kullback-Leibler (KL) distance. The KL distance has been defined and explained later. Thus during the process of adaptation of the bit stream to the available transmission rate the following three scenarios arise. We denote $r_i(t)$, $X_b(t)$, and $X_e(t)$ as the available transmission rate, the required base layer bit rate, and the enhanced layer bit rate respectively.

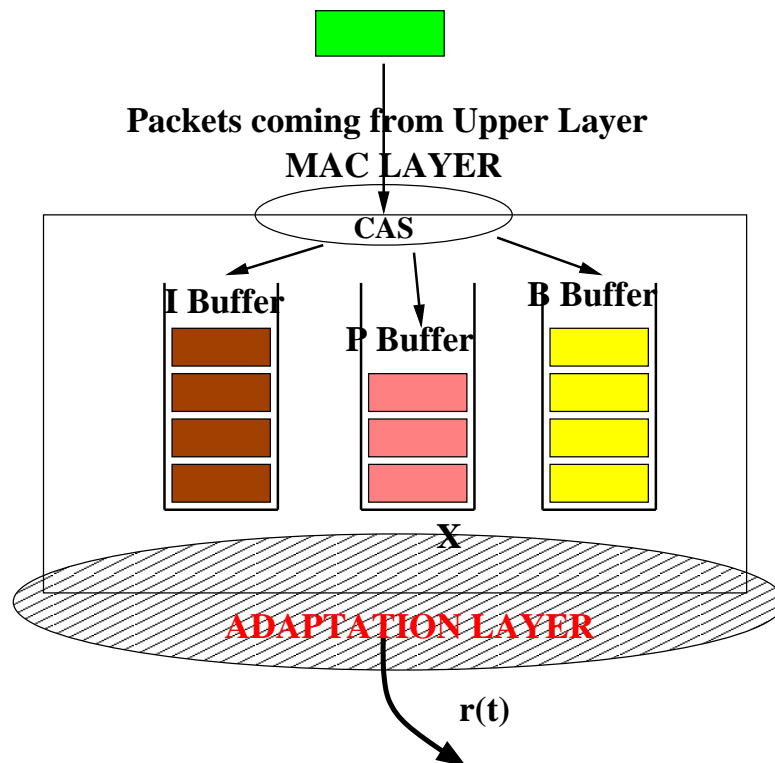


Figure 3.4. Adpatation Layer for Content Aware Rate Adaptation .

- *Case I:* $r_i(t) < X_b(t)$: The adaptation layer selectively drops base layer frames based on the proposed prioritization scheme.
- *Case II:* $X_b(t) \leq r_i(t) \leq X_b(t) + X_e(t)$: The highest priority layer is encoded. For the Enhanced Layer (EL), the adaptation layer encapsulates whatever portion of the EL that can be packed into the remaining available bandwidth.
- *Case III:* $r_i(t) > X_b(t) + X_e(t)$: This is the most favorable case where the encoding rate is higher than the bit stream rate and hence the entire data can be transmitted.

Through the proposed adaptation layer, the lowest possible granularity is achieved. However, a finite lag remains between the channel state and actual transmission of the data which in our case is the best achievable, i.e., lag by a slot. However, during packetization the granularity is measured in terms of the decision period. Depending on the current buffer status at both the transmitter and the user, and the current channel state, the rate control system needs to determine the encoding rate. The currently encoded packets would only be served after the existing MAC protocol data units (PDUs) are served. Thus rate control system need to determine the number of decision periods needed to serve the existing MAC PDUs. Hence, prediction of the rates by transition probability computation is important which we discuss next.

3.5.3 Transition Probability Calculation

We assume that the data rates are related to the distance of the user from the base station [17]. Though this assumption holds true under ideal conditions (i.e., no interference, no topological effects), we still use it for mathematical simplicity. If we assume that the system supports m rates, then a cell can be divided into m concentric rings; the innermost ring gets the maximum data rate and outermost ring gets the minimum. In one decision period, a user can go one ring above or below the current ring because of the limitation in speed. Let P_{uv} denote the probability of a user switching from a

data rate u to data rate v . Note, v can only be $u - 1$, u or $u + 1$. Let the steady state probability of a user being in the u th data rate be denoted by π_u . We define π_u as the ratio of the ring area by the cell area, and is given by $\pi_u = \frac{r_u^2 - r_{u-1}^2}{r_{m-1}^2}$ where r_u is the radius of the u th ring. Clearly, $\sum_{u=0}^{m-1} \pi_u = 1$. The transitional probabilities P_{uv} are computed geometrically (see appendix) and is as follows

$$\begin{aligned} P_{u \rightarrow u-1} &= \frac{2\sqrt{s(s-r_1)(s-r_2)(s-d)} - (\theta_1 r_1^2 - \theta_2 r_2^2)}{\pi r_2^2} \\ P_{u \rightarrow u} &= \frac{A_i}{\pi r_2^2} \\ P_{u \rightarrow u+1} &= \frac{A_{bge}}{\pi r_2^2} \end{aligned} \tag{3.16}$$

Note that these probability calculation will be different for the inner most and outer most rings.

The details of the above calculation are provided below. We are interested in calculating where the location of a user is currently in the i^{th} ring. The outer radius of i th ring is r_i and inner radius r_{i-1} . The maximum distance that the user can go within a decision cycle is d_{max} . A uniform velocity profile of the user is considered. Hence, then the user is equally likely to be anywhere which lies in three different rings- $i - 1$, i and $i + 1$.

Case I: Area for Δxyz (A_{xyz}) is given by

$$A_{xyz} = \sqrt{s(s-r_1)(s-r_2)(s-d)}$$

where s the sub-perimeter of the ΔXYZ is given by

$$s = \frac{r_1 + r_2 + d}{2}.$$

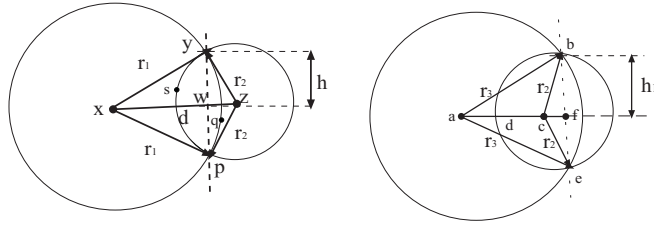


Figure 3.5. Proability Caculation for case I and II.

Let h be the height of the \triangle and it can be computed from the following relationship

$$h = \frac{2 \times \sqrt{s(s - r_1)(s - r_2)(s - d)}}{d}$$

On computing the height h which is the same for $\triangle XYW$, $\angle YXW$ (denoting $\angle YXW = \theta_1$) can be computed as $\theta_1 = \sin^{-1}\left(\frac{h}{r_1}\right)$. Thus the area of the arc XPY (A_{XPY}) is $\frac{\theta_1}{\pi} \times (\pi r_1^2)$. Similarly considering $\triangle ZPY$, $\angle YZW$ (let $\angle YZW = \theta_2$), the area of the arc AZYP can be computed as $A_{ZPY} = \frac{\theta_2}{\pi} \times (\pi r_2^2)$. Thus the area PQYS (A_{PQYS}) is

$$\begin{aligned} A_{PQYS} &= 2 \times A_{XYZ} - (A_{XPY} - A_{PQYS}) \\ &= 2 \times \sqrt{s(s - r_1)(s - r_2)(s - d)} - (\theta_1 r_1^2 - \theta_2 r_2^2) \end{aligned}$$

Case II: Similar to Case I, we proceed to calculate area of $\triangle abc$, height (h_1) of the triangle and $\angle bac$ illustrated in Figure 3.5 (let $\angle bac = \theta_3$ and $\angle bce = \theta_4$) as follows

$$\begin{aligned} A_{abc} &= \sqrt{s(s - r_3)(s - r_2)(s - d)} \\ h_1 &= \frac{2 \times \sqrt{s(s - r_3)(s - r_2)(s - d)}}{d} \\ \theta_3 &= \sin^{-1}\left(\frac{h_1}{r_2}\right) & \theta_4 &= \sin^{-1}\left(\frac{h_1}{r_3}\right) \end{aligned}$$

where s is given by $s = \frac{r_1+r_2+d}{2}$. Thus, the area of the arc abc is $A_{abke} = \frac{\theta_3}{\pi} \times (\pi r_2^2)$. Thus the reflex angle $\angle bce = \theta_5 = 2\pi - \theta_4$ and hence area of arc (A_{bde})

$$A_{bde} = \frac{\theta_5}{\pi} \times (\pi r_3^2)$$

Thus the desired area (A_{bge}) can be computed as follows

$$A_{bge} = \pi r_2^2 - ((A_{abke} - 2 \times A_{abc}) + A_{bde})$$

Thus the area of the middle strip is given by

$$A_i = \pi r_2^2 - A_{bge} - A_{pqys}$$

3.5.4 Analytical Modeling of Rate Adaptation

The encoding rate $r_i(t)$ for user i depends on the channel state, buffer status and bandwidth allocated as already discussed. This necessitates the computation of the number of decision periods needed to serve the existing MAC PDUs. The number of decision periods is basically the lag by which the encoder follows the channel state. If ξ is the number of decision periods needed to serve the MAC buffer B_i for user i , then

$$\sum_{l=0}^{\xi} S_i^l \times (P_{uv}^l \times R_i^l(t)) = f_b \times B_i \quad (3.17)$$

where $0 \leq f_b \leq 1$ is the buffer fullness, S_i^l is the number of slots allocated to user i in the l th cycle, P_{uv}^l is the transition probability from state u to v (already computed in Section 3.5.3, and $R_i^l(t)$ is the rate estimated using Equation 3.15- all in the l th decision cycle.

We need to compute the amount of data, $\chi_i(t)$, that the MAC scheduler would be able to support for the corresponding decision period. Depending on the scheduling paradigm (TM/FCTM/QFCTM) which determines the number of slots being allocated to the user, the amount of data which the adaptation layer may encode is given by

$$\chi_i(t) = S_i^\xi \times (P_{uv}^\xi \times R_i^\xi(t)) \quad (3.18)$$

Conversely, if the mobile terminal undergoes favorable channel condition (i.e., supports a higher data rate), then the rate controller would be tempted to transmit more data. Transmitting more data than required would compensate if the channel conditions deteriorate. The available buffer space ($v_i(t)$) during that decision period is

$$\begin{aligned} v_i(t) = & (1 - f_{rb}) \times RB_i - \sum_{l=0}^{\xi} S_i^l \times (P_{uv}^l \times R_i^l(t)) \\ & + \rho_i \times (\tau\xi) \end{aligned} \quad (3.19)$$

where $0 < f_{rb} < 1$ is the receiver buffer fullness, RB_i is the receiver buffer size, ρ_i is the rate of the playout curve rate, and τ is the time of each slot ($\tau = 1.67 \text{ ms}$ for HDR). The playout curve rate is the rate at which the player at the client end consumes the data.

The adaptation layer encoding rate considering the buffer constraints at both the transmitter and receiver sides is given by

$$r_i(t) = \min(\chi_i(t), v_i(t)) \quad (3.20)$$

Note that $r_i(t)$ is a random variable which gives the value of rate adaptation. Let $p(\cdot)$ denote the pdf for the random variable $r_i(t)$. Let $f_g(\cdot)$ denote the pdf of the Group of Picture (GOP) size distribution of the g th video. We use the Kullback Leibler (KL)

distance to characterize the performance of our adaptation scheme.

Definition 1: The Kullback Leibler distance gives the relative entropy between two distributions. This entropy is given as

$$D(p||q) = \sum_{x \in X} p(x) \ln \frac{p(x)}{q(x)} \quad (3.21)$$

Thus, the KL distance $D(p||q)$ measures the inefficiency by which the distribution $q(x)$ differs from distribution $p(x)$. Hence Equation (3.21) provides a metric to determine the lag or closeness of $q(x)$ to $p(x)$. Substituting $p(x) = p(\cdot)$ and $q(x) = f_m(\cdot)$, we define that the adaptation (\mathcal{A}) of the video to channel rate as

$$\mathcal{A}(p||f_m) = \sum p(\cdot) \ln \frac{p(\cdot)}{f_m(\cdot)} \quad (3.22)$$

The granularity of the rate of adaptation is bounded by the decision period as given in Equations (3.17), (3.18), (3.19), and (3.20) and hence we analyze the adaptation with respect to the decision period.

3.6 A Case Study: MPEG-4 over HDR

In order to verify the proposed rate adaptation schemes , we perform a small case study of MPEG-4 over HDR.

3.6.1 Content Aware Scheduling (CAS)

The MAC Scheduler is responsible for allocating slots among the users requiring different data rates. It can be any generic scheduler like the PF scheduler or any of those discussed in Section 3.4. These standard MAC schedulers do allocate slots either

based on the supportable channel rate and/or by the QoS demanded by each type of user but unfortunately is unaware of the content type being served from the MAC buffer. Hence, the MAC scheduler is incapable of exploiting the application specific features while scheduling. To overcome the shortcomings of the MAC scheduler we propose Content Aware Scheduler (CAS) as shown in Figure 3.6 which works in conjunction with the MAC.

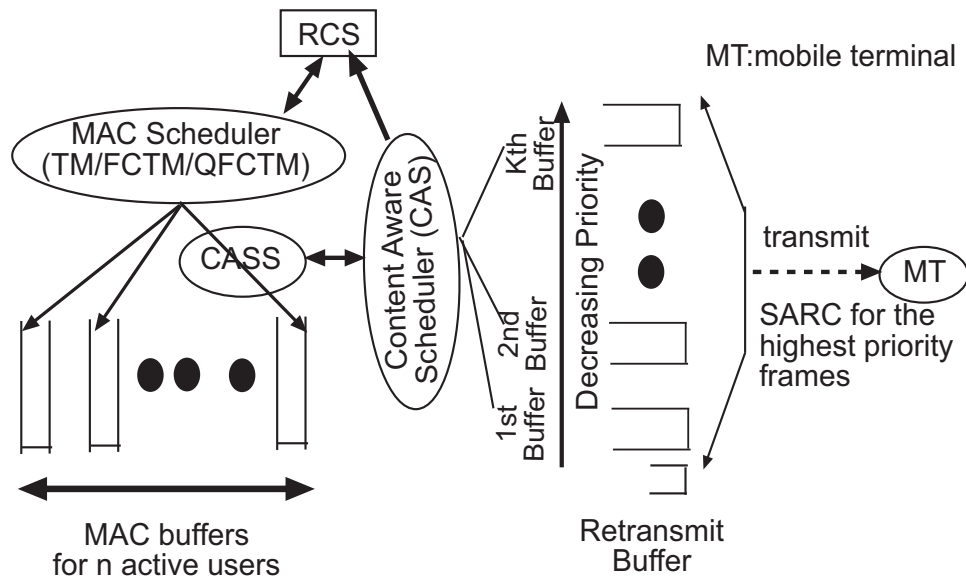


Figure 3.6. MAC Scheduler and CAS Scheduler.

Depending on the application, appropriate number of priority levels are generated. We assume that K can be any number of priority levels depending on the application. The CAS models the MAC buffer into K number of queues one for each priority level. For example in case of the FGS MPEG-4 video data there would be 4 queues one for each of the I, B, P-frames and the last one for the enhancement layer with the highest priority being given to the I-frames while the least to the enhancement layer. The generic MAC-S does not differentiate the priorities and simply serves the content in a first come first serve

(FCFS) queueing discipline. However, the proposed CAS being content aware enhances the system goodput. In many cases for real time streaming applications, transmission of the higher priority content is essential for the success of the application. In addition, the proposed CAS is aware that the packets need to be transmitted before a certain deadline at the receiver side. The CAS is smart enough to discard packets if *stale*. For i th packet, we define a boolean variable $SP(i)$ which determines whether the packet is stale or not. A packet i is defined stale if $SP(i) = 1$ and the condition is $t_{curr} > TS_i + t_{prop}$ otherwise $SP(i) = 0$ meaning the packet is not stale. Here t_{curr} , TS_i and t_{prop} are the current system time, the time stamp of packet i prior to which it needs to be transmitted, and the propagation time. This simple yet content aware scheduling mechanism prevents error propagation and increases the goodput.

The general MAC scheduler is unaware of the prioritization scheme and is unable to schedule efficiently. However, the CAS utilizes the priority information and also employs a Selective Adaptive Retransmit Control (*SARC*) mechanism for the highest priority packets. The SARC mechanism can be utilized till any level of priority (i.e., for each buffer) but for the sake of simplicity, we restrict ourselves to the highest priority buffer. Additionally, the CAS employs *Fast Transmit Scheme (FTS)* to take advantage of favorable channel condition but does take into account of the upper limit In summary, the rationale for employing SARC is as follows :

1. Provide unequal error protection (UEP) to high priority data and thereby increase the transmission quality and prevent error propagation.
2. Transmit the mobile terminal's buffer fullness (f_{rb}) back to the base station through the SARC mechanism. f_{rb} is essential for determining boundary condition as discussed in Section 3.5.

The SARC mechanism of CAS keeps retransmitting packets from queue with the highest priority till the mobile terminal acknowledges the successful arrival of the frame.

The number of retransmissions is bounded by the timing requirements, i.e., till the packet becomes stale. Continuing with the FGS MPEG-4 example the I-frame being the most important video data frame, multiple retransmissions might be required for it. Missing or corrupted I-frame results in wastage of the corresponding B, P and enhanced layer frames even if they arrive correctly. The FTS scheme enables the CAS to transmit more data than the playout curve rate of the mobile terminal when the channel conditions are favorable provided CAS has assigned sufficient slots by the MAC-S. Since the receiver buffer of the mobile terminal is finite, the CAS should restrain from transmitting data that would overflow the mobile terminal buffer. Similar to the ARQ mechanism in TCP, the mobile terminal transmits the available buffer space in each acknowledgment packet. The CAS makes the rate control system module aware of β . Thereafter, RCS module takes into account of the available rate, buffer space at both the MAC and mobile terminal, and computes the upper bound for data transmission as derived in Section 3.5. It also determines whether the fast transmission scheme outlined earlier is achievable or not. The working principle of CAS is explained below. Without loss of generality let us explain the CAS scheduling of an application having k priority levels and the highest priority level enabled with SARC in the K^{th} decision cycle of the MAC scheduler.

Algorithm 6 Content Aware Scheduler

```

1: CAS  $\leftarrow$  control state
2: if (Overflow Constraint not violated) then
3:   if 1HP not empty then
4:     if  $TS_{1HP\text{-packet}} < TS_{2nd,\dots,kth\text{-packet}}$  && !SP(1HP) then
5:       CAS schedules 1HP-packets
6:       CAS employs SARC
7:       //wait for  $\zeta$  slots for ACK
8:     end if
9:   end if
10:  if CAS  $\leftarrow$  SARC feedback then
11:    Retrieve retransmit packet numbers
12:    Compute Overflow Constraint
13:    if !SP(1HP-packet) then
14:      CAS retransmits 1HP-packets
15:      CAS employs SARC
16:      if slot available then
17:        pack  $2_{nd}HP, \dots, k_{th}HP$ -packets respectively
18:        //wait for  $\zeta$  slots for ACK
19:      end if
20:    end if
21:  end if
22:   $i=2$ 
23:  if  $i$ th-Buffer not empty then
24:    CAS transmits  $i_{th}$  Buffer packets.
25:    Increment  $i$ 
26:  end if
27: end if

```

We present the CAS in Figure 6 where $1HP$, TS and SP respectively denote packet from the highest priority queue, timestamp of the packet and whether the packet is stale or not. The CAS initiates transmission with the highest priority packets employing SARC and following up with the rest of the frames without employing any ARQ on the rest of the packets. For CAS, goodput would be more effective measure than throughput. We define *goodput* for real time data as the number of packets transmitted per decision cycle by the CAS scheduler that the mobile terminal successfully utilizes. Next we analyze the proposed CAS algorithm.

3.6.2 CAS Goodput

The analytical modeling of the CAS and the goodput of the system using CAS are presented. First, we model the SARC mechanism. Let $P_b(m)$ and $P_e(m)$ denote the probability of the bit error and packet error respectively for sending m packets simultaneously. If $P_s(m)$ denotes the probability of successful packet transmission then

$$P_s(m) = 1 - P_e(m) = [1 - P_b(m)]^L \quad (3.23)$$

where L denotes the length of the MAC protocol data unit (PDU). To provide more protection to the highest priority packets, error correction code (ECC) is employed. The modified probability for packet transmission with ECC parameter as t_{ecc} is given by

$$P_s(m, y) = \sum_{i=0}^y \binom{L}{o} P_b(m)^i [1 - P_b(m)]^{L-o} \quad (3.24)$$

where o is the number of errors corrected and y is the total number of correctable bit errors. Thus $P_s(m, t_{ecc})$ denotes the probability for successful packet transmission for highest priority packets. However, the number of retransmissions for these packets is

limited by ϵ , the maximum number of retransmits for the i th 1HP packets. We compute ϵ as

$$\epsilon = \frac{TS_{1HP\text{-packet}}}{t_{prop} + t_{arq}} \text{ if !SP(1HP-packet)} \quad (3.25)$$

where t_{prop} , t_{arq} and TS denote the propagation, the time after which the ACK transmitted by the mobile terminal is received, and the timestamp of the corresponding packet before being marked stale, respectively. The mean of the number of retransmissions for 1HP-packets is given by

$$\delta_{1HP} = \sum_{i=1}^{\epsilon} (1 - P_s)^{(i-1)} P_s \times i = P_s \times \frac{1 - P_s^{\epsilon}}{(1 - P_s)^2} \quad (3.26)$$

Let n_k be the total number of MAC-PDUs generated by FPPL for the k th data segment and let n_{k1HP} and n_{kR} be the number of 1HP-frame packets and $2, \dots, k$ - frame packets respectively such that $n_k = n_{k1HP} + n_{kR}$. The timestamp of the frames is dependent on the data size distribution. In Section 3.5, we have already defined $f_q(\cdot)$ as the pdf of the data size distribution and $f_q^j(t)$ as the pdf of the timestamp of the j th data of the q th real time stream. We model δ_{Rj} as the parameter which determines whether $2, \dots, k$ - frame are stale or not consider that $\delta_{Rj} = SP(j)$. Thus the effective number of packets transmitted per decision period of CAS is given by

$$\Delta_i = n_{kI} \times \delta_I + n_{kR} \times \delta_{Rj} \quad (3.27)$$

Therefore, the average goodput (ρ_{gp}) of CAS is given by

$$\rho_{gp} = \frac{\sum_{i=1}^K \Delta_i}{K} \quad (3.28)$$

where K is the total number of decision periods for which the video is transmitted. Note that throughput of such systems would fail to capture the actual system performance since it would not consider adaptive selective retransmission and selective frame dropping based on packet staleness. In Figure 3.7, we compare the CAS goodput to QFCTM throughput.

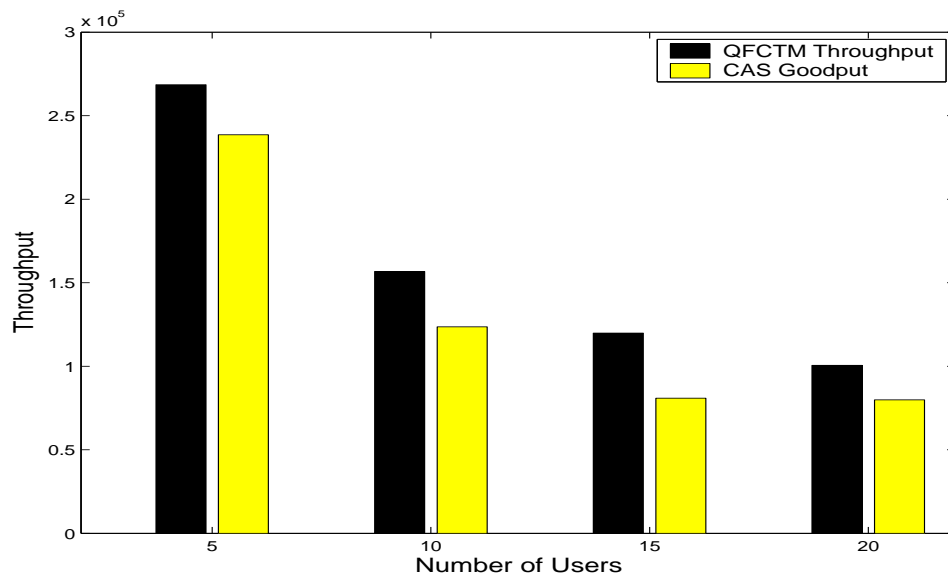


Figure 3.7. Comparison of CAS Goodput and QFCTM Throughput.

The CAS goodput trails the QFCTM throughput. In an ideal scenario, with no packet loss and zero lag, the CAS goodput should have been equal to the QFCTM throughput. Also, note that with increase in number of users, the CAS goodput is penalized due to QoS constraints.

3.7 Simulation and Experimental Results

We conducted simulations where a single multi-rate cell with multiple users were simulated to demonstrate the performance of the proposed estimation and adaptation technique with respect to the existing scheme for HDR. The users were randomly dis-

tributed over the cell. All our results hold for vehicular users since the speed of the users vary between 9 to 15 meter/second. The data rates used were standard HDR data rates [13] 38.4 Kbps to 2.4 Mbps. Random mobility of users were assumed and the data rates were updated according to the distance from the base station. To evaluate our proposed scheme, we focussed on two aspects : **(i)** the accuracy of the prediction scheme and **(ii)** the performance of the proposed scheduling algorithms in terms of the number of satisfied users and throughput achieved.

3.7.1 Simulation Results

We present simulation results to illustrate the performance of the proposed estimation technique with respect to the existing PF scheme for HDR. However, please note that similar results would also hold true for existing multi-rate wireless systems. In Figure 3.8, the actual channel rate, the estimated channel rate for both the current and proposed estimation techniques are presented. The proposed scheme outperforms the existing PF scheme since the adaptation for the proposed scheme takes place on per slot basis. In Figure 3.9, we magnify the time block (400 to 600 slots) of Figure 3.8 to show that the proposed technique adapts to the time varying channel conditions within a decision period whereas the existing estimation technique remains static during each decision period. We further present estimation error for the proposed and existing scheme in Figure 3.10. Next, we present the scheduling algorithms based on the estimates of the data rates.

In Figure 3.11, we present the system throughput achieved for the round robin (RR), TM , FCTM and QFCTM scheduling algorithms. As expected, the TM algorithm achieves the highest throughput. Noted that all these algorithms achieve these throughputs based on the proposed channel rate estimation techniques. The system throughput

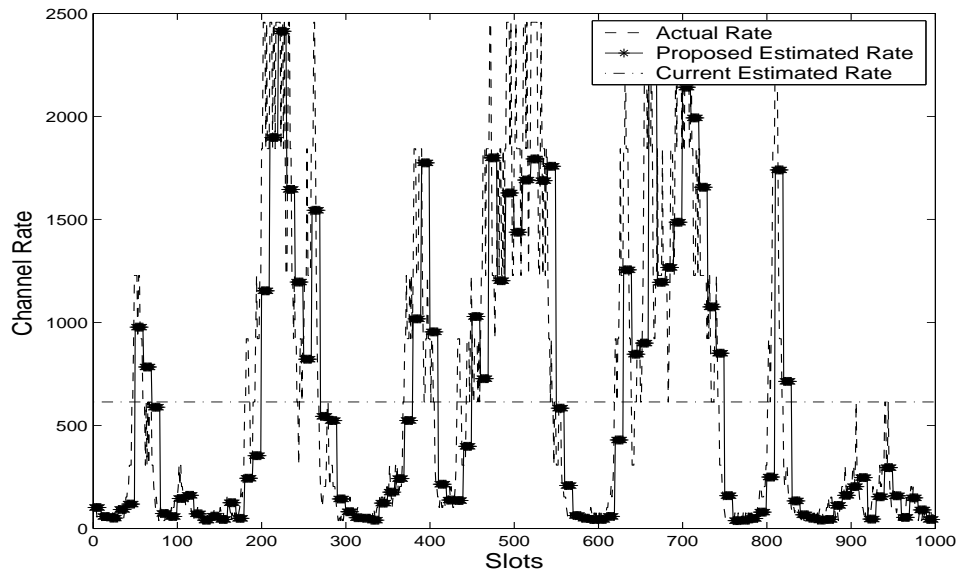


Figure 3.8. Estimation and adaptation for proposed and current Scheme.

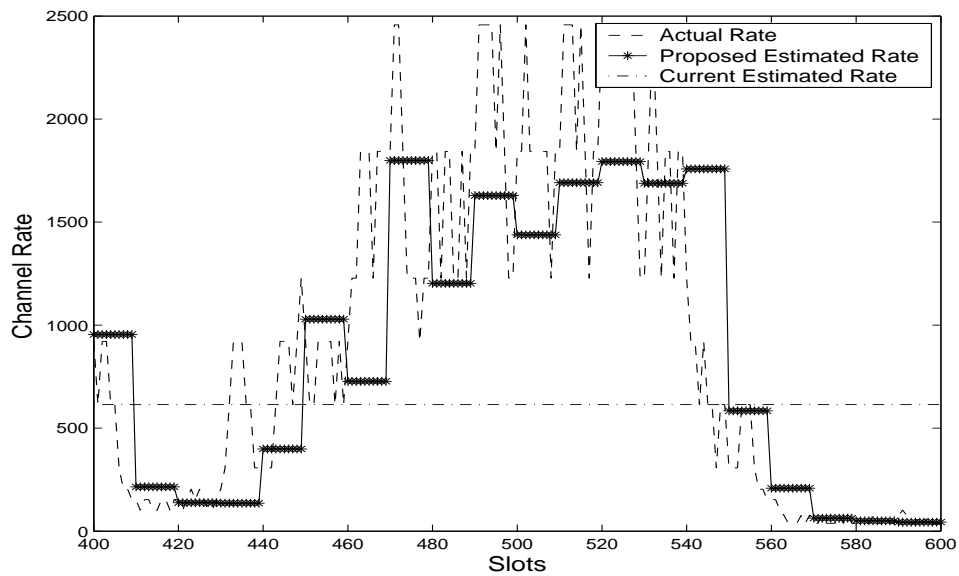


Figure 3.9. Magnified view (Time epoch 400 to 600).

is measured with respect to the number of users in the system to conclude whether multi-user diversity gain is being exploited or not.

The results confirm with the notion that higher the fairness and QoS constraints lesser is the throughput achieved. The per user throughput for different scheduling

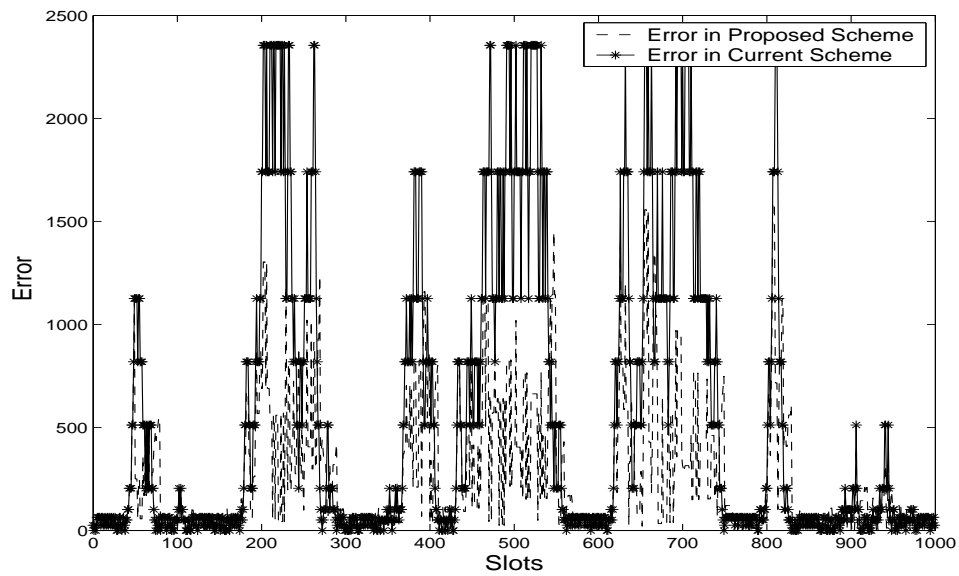


Figure 3.10. Error in the Proposed and PF Scheme.

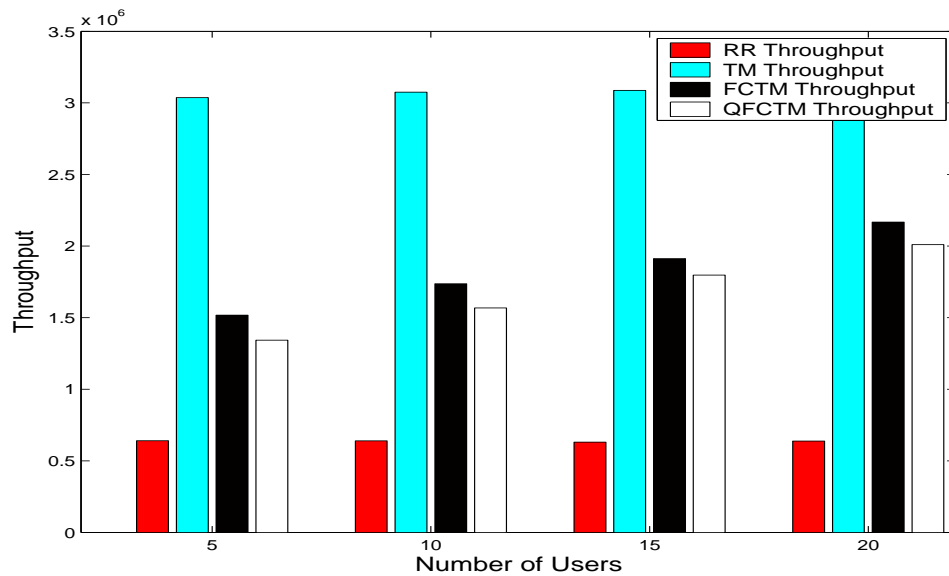


Figure 3.11. System Throughput for different scheduling algorithms.

algorithms is presented in Figure 3.12. However, per user throughput would not reflect the fairness of the scheduling disciplines. To determine the fairness of the schemes, Figure 3.13 presents the standard deviation (SD) for the per user throughput. The results show that QFCTM has the least SD and hence is the most fair. Thus, we have successfully

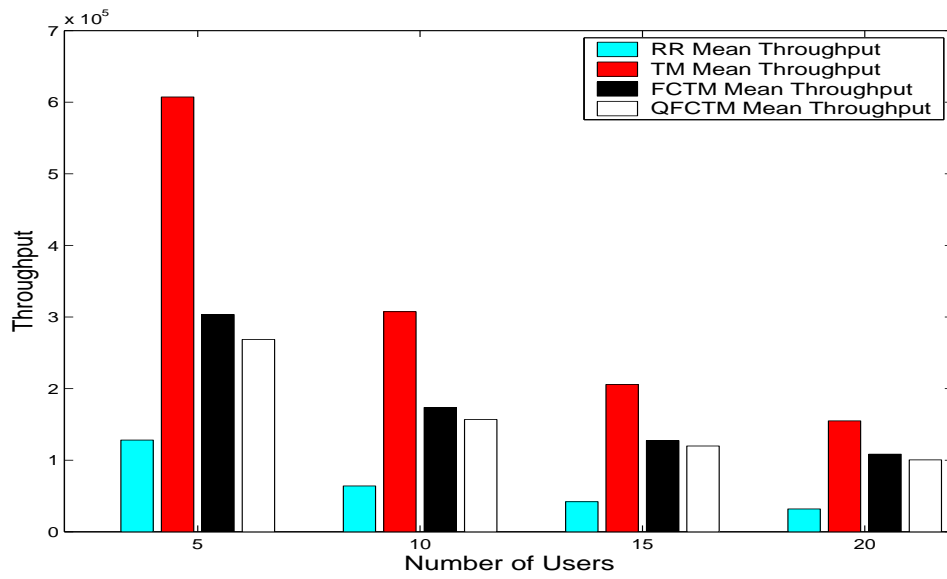


Figure 3.12. Mean throughput per user for different algorithms.

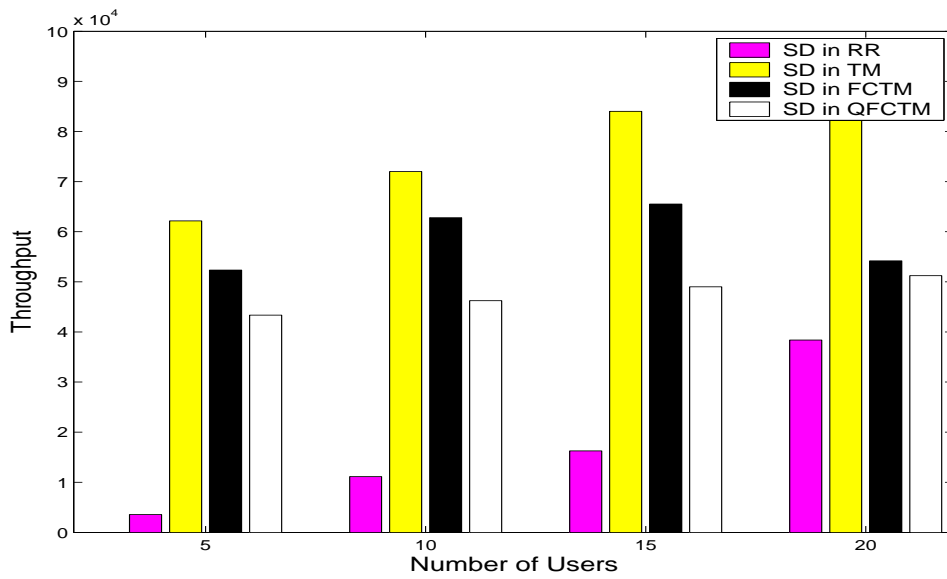


Figure 3.13. Standard deviation (SD) of throughput for different algorithms.

adapted to the channel rate and provided scheduling algorithms which guarantee fairness, QoS constraints and maximizes throughput under the above constraints.

In Figure 3.14, we evaluate the variation of the different algorithms in terms of percentage of the satisfied users with the variation in schedule length for our proposed

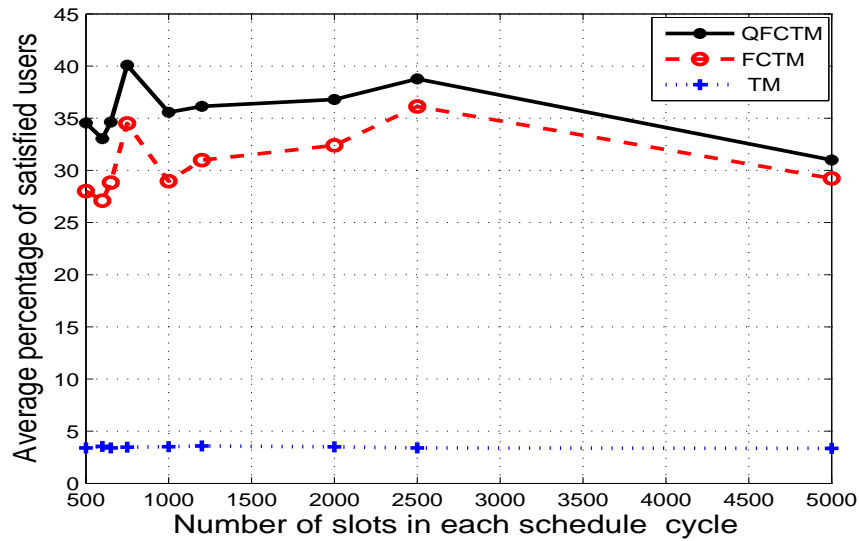


Figure 3.14. Variation on the Percentage of Satisfied Users for different algorithms with respect to number of slots in the schedule cycle.

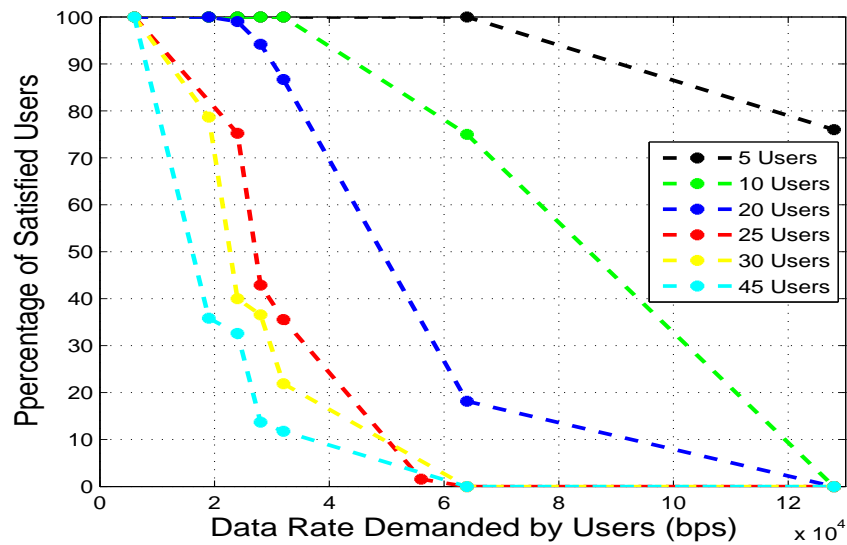


Figure 3.15. Variation of the Percentage of Satisfied Users for QFCTM for different data rates.

estimation scheme. From our simulation results, we recommend that with fairly accurate channel rate estimation scheme like the proposed one, a schedule length between 800 to 1000 slots would be ideal for maximizing the number of satisfied users. We also focussed

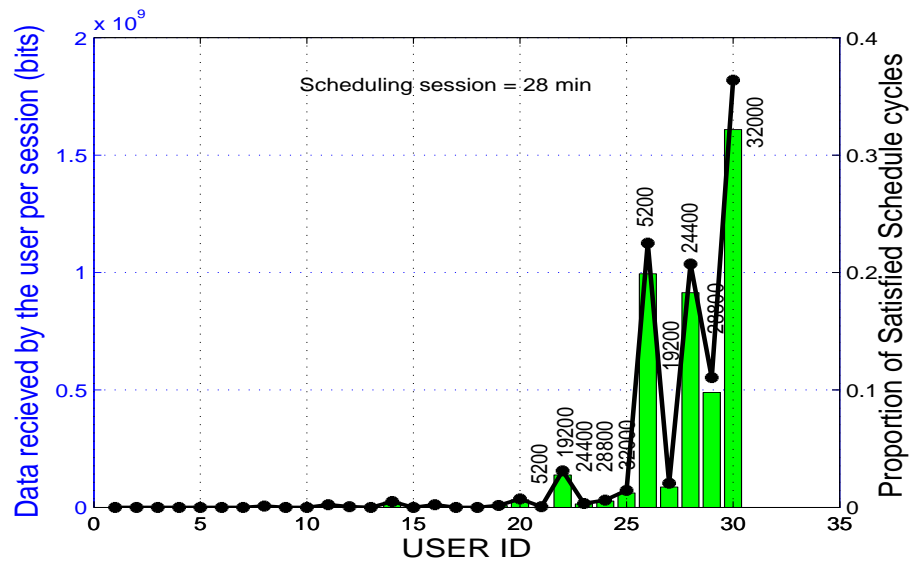


Figure 3.16. Fairness for TM algorithm.

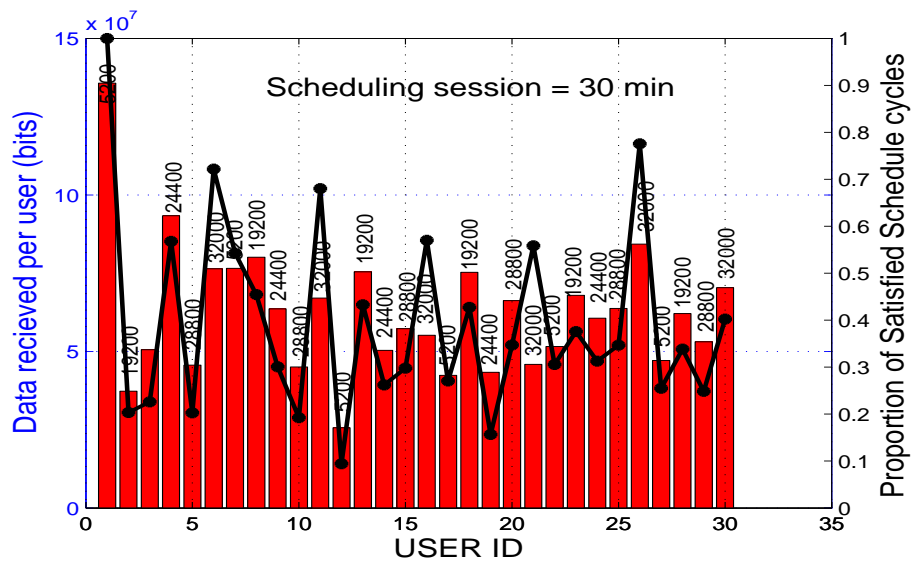


Figure 3.17. Fairness for FCTM algorithm.

on finding out the cutoff values for the number of satisfied users with specific data rates. Figure 3.15, shows the variation of number of satisfied user for different number of total users for QFCTM. Here for each case, the data rates for all users were set to the same value.

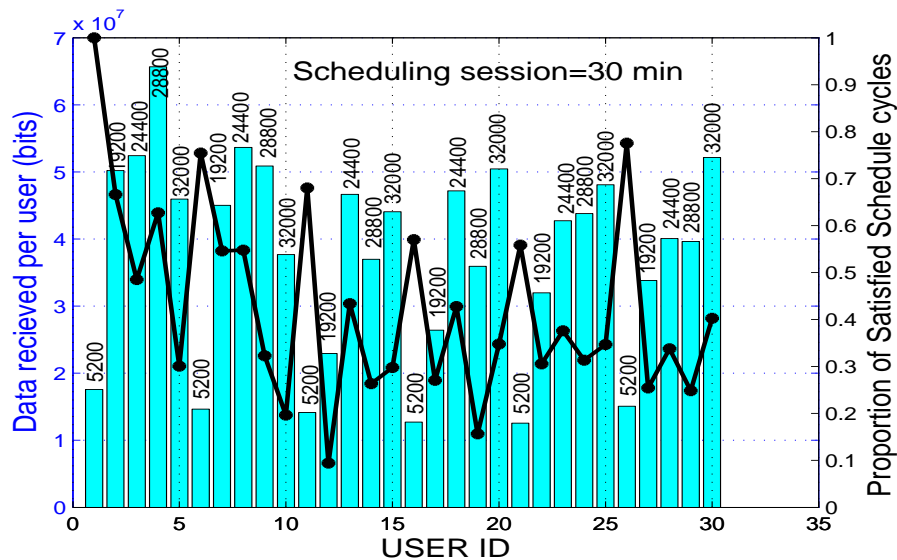


Figure 3.18. Fairness for QFCTM algorithm.

In order to demonstrate the fairness aspect of each algorithm, a pool of 30 users were chosen. Each of the m users had the same data requirements but suffered different channel state variations. In Figure 3.16, we note that most of the users are starved and only the users with good channel states have been served leading to a very low number of users getting satisfied. Both the FCTM and QFCTM algorithms support a larger number of satisfied users as shown in Figures 3.17 and 3.18. However, QFCTM supports a higher number of average satisfied users.

Resolution	Clip name	fps	bitrate
QCIF	paris	15fps	64kbps
QCIF	foreman	25fps	64 kbps
CIF	football	15fps	1Mbps

Table 3.1. Specification of the files used in our simulation.

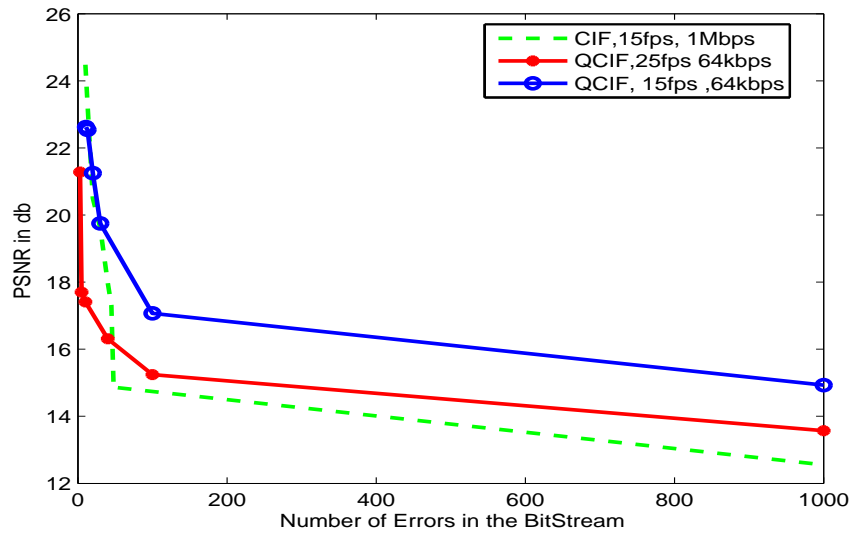


Figure 3.19. Variation of PSNR with BER using proposed CAS algorithm.

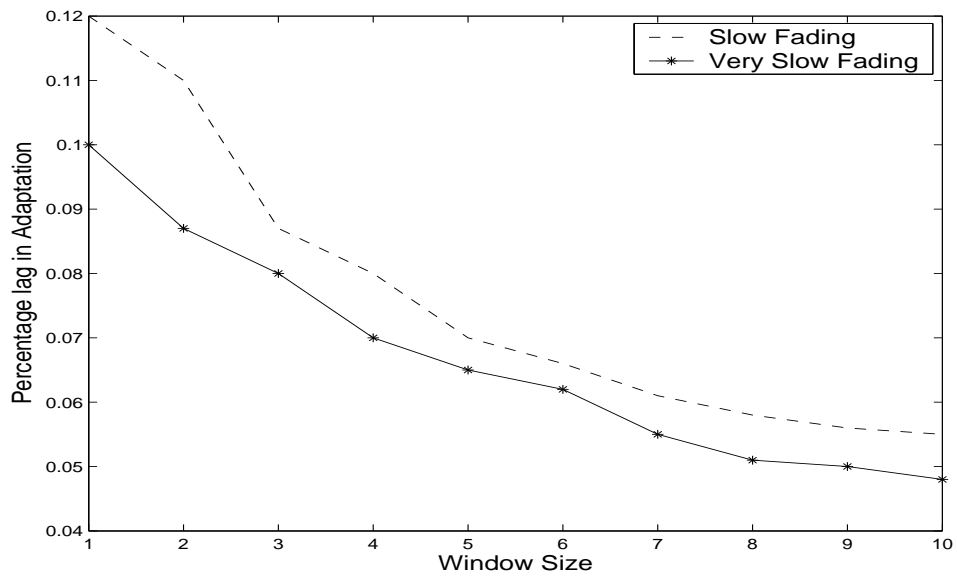


Figure 3.20. The variation of CAS adaptation due to mobility.

3.7.2 Results for Rate Adaptation and CAS

We have conducted simulation experiments to illustrate the performance of the proposed adaptation technique as well as the effectiveness of CAS with respect to the existing scheme for HDR. Figure 3.19 highlights the graceful degradation of Peak signal-

to-noise ratio (PSNR) values using the CAS scheduler. The video test sequences chosen comprises of both the simple profile (SP) and advanced simple profile (ASP). In order to ensure the performance of the proposed schemes, we use representative test sequences of *foreman*, *paris* and *football* which have varying resolution (Common Intermediate Format (CIF), Quarter CIF (QCIF)), frame rates and bit rates. The specifications of the streams are listed in Table 3.1.

As for the lag in adaptation, we analyzed equation 3.22 numerically. Figure 3.20 shows how the system is able to learn and adapt if a sufficiently long window of observation is allowed. As we do not deal with fast fading channels, we compared ‘slow’ and ‘very slow’ fading channels. Obviously, the adaptation is better for slower fading channels. In Figures 3.21 and 3.22, we show how the throughput of I-frames in multimedia streaming could be improved in the presence of CAS.

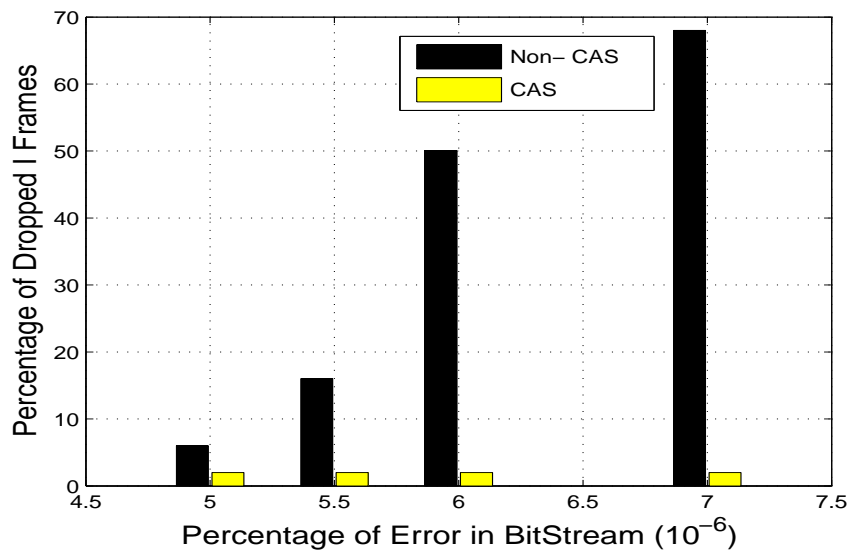


Figure 3.21. Improved throughput of I-frames due to CAS..

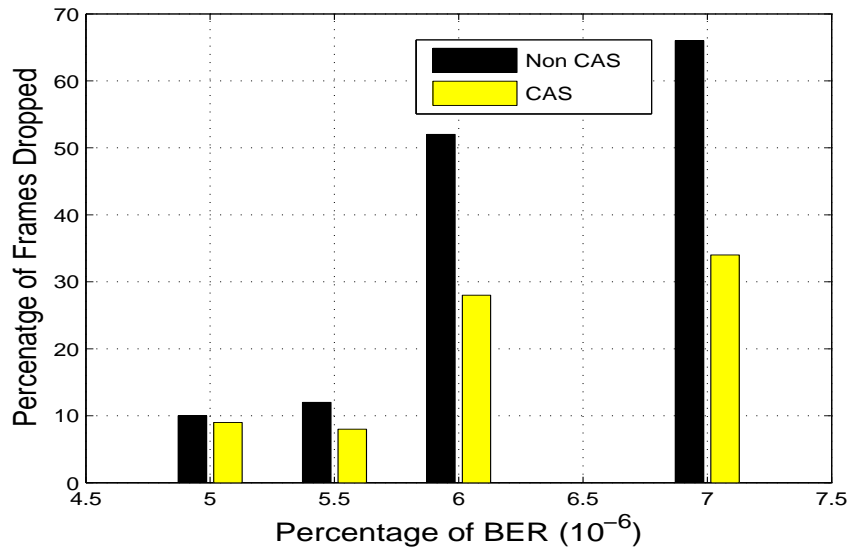


Figure 3.22. Error resiliency of CAS and throughput of I-frames..

3.8 Summary

In this chapter, we developed focussed on mechanisms to harness the higher bandwidth offered by multi-rate wireless systems. We emphasize that not only the variation of the data rates need to be tackled but also the varying demands of the different data services need to be considered for success of the multi-rate wireless systems. We have also shown that maximizing the number of satisfied users is not synonymous to maximizing the system throughput. We have utilized information theoretic metrics to estimate the channel state variations. Specifically, a non-parametric estimator of Renyi's entropy using the Parzen windowing technique has been employed to estimate the probability density function of the channel rate variation. Based on the estimated channel rates, scheduling algorithms which guarantee GPS fairness, and assure data rate and QoS requirements have been developed. Simulation, experiments and numerical analysis that our proposed mechanisms perform favorably.

CHAPTER 4

THEORETICAL BOUNDS FOR MAXIMUM NUMBER OF SATISFIED USERS IN MULTI-RATE WIRELESS SYSTEMS

The insight of the maximum possible number of satisfied users that can be supported in multi-rate wireless system has not been investigated. Existing opportunistic scheduling algorithms are effective in exploiting channel variations and maximizing system throughput in multi-rate wireless networks. Most scheduling algorithms ignore the per-user quality of service (QoS) requirements and try to allocate resources (e.g., the time slots) among multiple users. This leads to a phenomenon commonly referred to as the *exposure problem* wherein the algorithms fail to satisfy the minimum slot requirements of the users due to substitutability and complementarity requirements of user slots. Thus, existing schemes ignore maximizing the number of satisfied users. To eliminate this exposure problem, in this chapter we propose a novel scheduling algorithm based on two-phase combinatorial reverse auction with the primary objective to maximize the number of satisfied users in the system.

The rest of the chapter is organized as follows. In Section 4.1, we formulate the scheduling problem for a multi-rate, multi-user time-slotted system. Section 4.2 models the user utility, identifies the exposure problem and studies its implications on the user utility. The mapping of combinatorial auction to multi-rate scheduling is presented in Section 4.3. A scheduling scheme based on reverse combinatorial auction is presented in Section 4.4. We analyze the proposed algorithm and compare its performance with other existing schemes in Section 4.5. Simulation results are presented in Section 4.6 followed

4.1 Problem Formulation

Time constraint scheduling is a necessity for delay-sensitive applications. We justify the above by explaining the timing requirements of VoIP applications. According to the

International Telecommunication Union (ITU-T) G.114 specifications [33], for good and pleasing voice quality, the end-to-end delay for both the forward and reverse path should not be more than 150 ms. This delay is contributed by various sources: (i) the voice coder with a processing delay of 10 ms, (ii) bit compression module with delay upto 7.5 ms, (iii) packetization scheme that introduces a delay between 20 to 60 ms, (iv) serialization with varying delay between 0.20 ms to 15 ms, (v) queuing/buffering and network switching delay of around 65 ms, and (vi) the de-jitter buffer with worst-case delay figure of 40 ms. Summing up these figures, it is easy to observe that the delay budget is already in excess of the acceptable ITU-G.114 requirements. That too is without taking into account the last hop wireless link where additional delay may occur due to uncertainty associated with the underlying wireless channel. Thus, to keep the end-to-end delay within acceptable limits, the wireless delivery system must schedule user data delivery within strict timing constraint. Therefore, the objective of scheduling is not only to improve the *throughput* of the system and enforce fairness among participating users, but also to meet the *minimum data requirements* of users at each scheduling time slot. It is not possible to provide such delay-sensitive scheduling with the help of existing scheduling techniques. We have developed techniques to evaluate the maximum number of users that can be satisfied in such multi-rate wireless systems for real-time data requirements.

4.1.1 Contributions

We take a fresh approach to the delay-sensitive scheduling problem by borrowing techniques from auction theory [37] and strive to find the maximum possible number of satisfied users. We consider a cellular network with one base station and multiple users. The resources available to the base station (e.g., time slots, frequency bands, codes) form the goods which are sold to the users in a *market*-like environment. The users value these goods distinctively and express the values in terms of a common transaction unit called

money[†]. We also consider a time-slotted wireless packet data system where the duration of an individual time slot is smaller than the average fading duration of the received signal. Thus, during symbol transmission, we can assume that the underlying wireless channel exhibit time-invariant properties.

Each user demands a certain number of slots (called bundle and denoted as S_b) in order to satisfy the *minimum* data requirement within a specific schedule cycle. The number of such slots depends on the condition of the underlying wireless channel. Since the market has multiple indivisible goods and each user's individual valuation of the goods depends on the *bundle* of goods received, we formulate the scheduling problem as a specific case of *combinatorial auction*. This is due to the fact that a single item transaction of the goods do not suffice since the user is more interested in the *sum total of the data received*. This underlying condition is exactly the reason why *single* slot allocation approach is not appropriate for delay-sensitive applications in multi-rate wireless systems. Consequently, the schedulers based on the principles of opportunistic scheduling are unable to satisfy the minimum data rate constraints demanded by the users.

In contrast, scheduling based on combinatorial auction deals with multi-slot allocation. Our proposed scheme can be used to satisfy the minimum data rate constraint of individual users. To model our system, we use both forms of combinatorial auction – *forward* and *reverse*. In the forward auction there exists a single seller who wants to sell multiple distinct goods to multiple buyers, while in the reverse auction there exists a single buyer who wants to procure goods from multiple sellers. In the former case, the intention of the seller is to maximize the total money received, whereas in the later the buyer tries to choose from sellers who quote the minimum price.

In our study, we establish that existing opportunistic scheduling algorithms are

[†]We use the concept of "money" as a tool for defining the resource allocation problem and as such has no significance in real life.

at best equivalent to our proposed scheme. We first formulate a combinatorial forward auction based multiple slot scheduling scheme that guarantees the minimum data requirements of the users. However, such an approach is shown to be \mathcal{NP} -complete [37] and hence computationally intractable. To design a tractable solution, we therefore re-formulate the problem based on the reverse auction and propose an approximation algorithm.

The main contributions are summarized as follows:

- We demonstrate that most of the existing scheduling algorithms suffer from the *exposure problem* and hence fail to guarantee the minimum data requirements of the admitted users.
- We use *combinatorial reverse auction* to formulate the scheduling problem with two different objectives: (i) guarantee the minimum data rate of the users, and (ii) maximize the overall system throughput.
- By mathematical analysis, we show that the proposed scheme is capable of supporting maximum possible satisfied users with hard real time requirements than the existing schedulers. Our approach also leads to significant gain in the system throughput.
- We prove that the worst case performance of the proposed approximate algorithm is bounded by a multiplicative factor $(1 + \log m)$ corresponding to the optimal solution, where m denotes the number of slots in a schedule cycle. We have also derived the time complexity of the algorithm.
- We conduct simulation experiments to evaluate the performance of our proposed algorithm with respect to two extreme scheduling disciplines: round-robin and throughput maximization. It is observed that our approach can schedule more users whose minimum QoS requirements are met than existing schemes.

- Finally, we propose a design parameter, α , that determines the trade-off between guaranteeing the user utility level (a measure for user satisfaction) and system throughput. The variation of system capacity with the number of satisfied users for different scheduling algorithms is also shown.

We first describe the system model under consideration and qualitatively formulate the scheduling problem. We also define the objective functions for optimal scheduling. For the sake of completeness, we start by briefly describing the basics of auction theory which forms the basis of our proposed scheduling scheme.

4.1.2 Preliminaries on Auction Theory

An *auction* is the process of buying and selling goods by offering them up for bid (i.e., an offered price), accepting bids, and then selling the item to the highest bidder [37]. In economics, an auction is a method to determine the value of a commodity that has an *undetermined* or *variable* price. In some cases, there is a minimum or reserve price; and if the bidding does not reach the minimum price, no transaction between buyers and sellers is executed. Most of the auctions are primarily *forward auctions* which involve a single seller and multiple buyers. The buyers compete among themselves in order to procure the goods of their choice by placing an initial bid that they feel is an appropriate price for the item under consideration. However, in *reverse auctions*, the role of the buyers and seller are reversed. A buyer places a request to purchase a particular item and multiple sellers bid to sell the requested item. The winner of a reverse auction is the seller who offers the lowest price. Sometimes, the bidders are interested in bidding for *multiple items* at the same time. In such a combinatorial bid, the bidder offers a price for the *collection of goods* according to the choice of the bidder rather than placing a bid on each individual items separately. This results in *combinatorial auction* where the

the auctioneer selects a *set* of combinatorial bids that provides the maximum return in revenue without assigning any item to more than one bidder.

4.1.3 Scheduling in Wireless Networks: A Qualitative Formulation

Wireless users derive utility from the services received from the wireless service providers. The utility perceived is a function of the amount of data received in a *specific time epoch*. In our study, we define a non-zero minimum utility, U_{min} , that must be met for user satisfaction. Corresponding to U_{min} , there exists a certain *minimum* amount of data, D_{min} , that must be made available to each user within a specific deadline. Failure to transmit the “entire” D_{min} to the user within the deadline or the schedule epoch, results in two-fold penalty that not only leaves the user dissatisfied but also penalizes the system throughput since partial transmission of the data ($< D_{min}$) does not contribute towards increasing the user utility. A representative example is the scenario of streaming multimedia (MPEG-4 video) where delayed transmission of packets associated with any I-frame results in the frame being discarded [28]. Thus, in real-time scheduling systems such as multi-rate wireless packet networks, each user is assumed to require at least D_{min} bytes of data every schedule cycle. Hence, instead of solely maximizing the system throughput, our proposed scheduler aims at maximizing *the number of users whose minimum utility is guaranteed*.

If the minimum utility of a user cannot be satisfied within the schedule cycle, the scheduler does not grant any slots to the user to avoid in the two-fold penalty as discussed above. However, once the number of allocated users for the particular schedule cycle has been decided, the schedule then endeavors to maximize the utility among those users. In general, we argue that the throughput maximization assuming “pay-per-byte” philosophy is detrimental for maximizing the revenue of the service provider since it does not maximize the number of satisfied users whose minimum data rate (D_{min}) is

Table 4.1. Notations Used in auction based scheduling

Notation	Meaning
m	The number of slots that defines the schedule cycle
M	The set of slots available for auction
n	The number of users admitted by the system
N	The set of satisfied users
r_{ij}	Possible transmission bit rate for slot i supported by user j
w_{ij}	Schedule vector
D_{min}	The min. demand data rate per user per cycle
D_{max}	The max. data rate for which marginal utility increase is nonrecognizable
U_{min}	The min. utility derived by the user
U_{max}	The max. utility derived by the user
T_p	System throughput for a schedule cycle
Θ_P	Penalty function for system throughput loss per cycle
Θ_U	Utility derived by the users
$V()$	Mapping between the data received to the corresponding utility $U()$
\mathcal{B}_I	The initial feasible bid set for the forward auction
S_b	The number or bundle of slots available to a user
$P_j(S_b)$	The deprivation function for the bundle S
p_j	The price user j quotes for the bundle S
A_j	The set of slots acquired by user j
$f_j(S_b)$	Mapping function between the price and the deprivation function
\mathcal{A}^k	Initial feasible bundle or set for round k in the restricted phase
\mathcal{A}^l	The set of all remaining slots after round l of the restricted phase

guaranteed. It is thus rational to assume that the generated revenue is proportional to the number of users who are satisfied in the *long run* if the service provider wants to keep the churn rate (measure of the user attrition rate) under control [45].

4.1.4 Wireless System Model

We consider a single cell, multi-rate time division multiple access (TDMA) wireless data system supporting n users[†]. *Downlink* scheduling of the wireless frames is realized by the base station in a time division manner whereby in each time slot the data is transmitted to only one user, as in HDR based systems [13]. The schedule cycle, the rate supported by each user and the slot allocation for the multi-rate wireless system is

[†]In this study, we assume that the n users have already been admitted by the session admission control algorithm, specific details of which is beyond the scope of this paper.

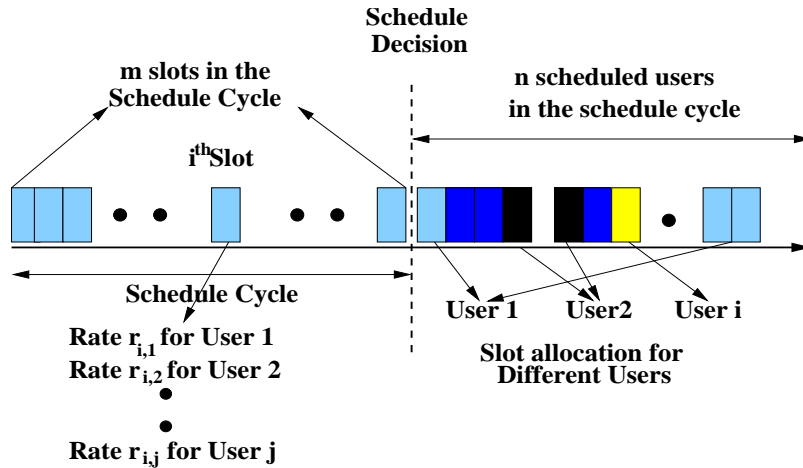


Figure 4.1. Illustration of the Schedule Cycle for Multi-Rate Wireless System.

illustrated in Figure 4.1. Table 4.1 lists the various parameters for system description and used by our proposed algorithm in Section 4.4. Through channel state prediction and feedback mechanisms, the base station is made aware of the channel quality and the corresponding data rate experienced by each user, for a specific time window corresponding to the schedule cycle. For each user, the slots in the schedule cycle comprise of the schedule vector.

Among the admitted users, let r_{ij} denote the possible transmission bit rates for slot i and experienced by user j . Consequently, $r_{i,j} \in \{0, r_1, \dots, r_R\}$ where R denotes the total number of feasible transmission bit rates and 0 signifies that the user is not allocated any slot in the schedule vector. Scheduling decisions are *periodic* and made every m slots (the actual value of m is implementation specific). We also denote the length of each slot as t_s ms. Thus, if a schedule decision is performed at time instance $T_d = a$, subsequent decisions are made at times $T_d = a + i \times m \times t_s$ for $i \geq 1$. Associated with every scheduling decision is the schedule matrix, $[w_{i,j}]$ where $1 \leq i \leq m$ and $1 \leq j \leq n$. If user j is granted slot i then $w_{i,j} = 1$, otherwise $w_{i,j} = 0$. We also

introduce a value function, $V()$, that maps the data received in a particular slot to the corresponding satisfaction or utility, $U()$, of the user receiving in that time slot.

4.1.5 Optimal Scheduling of Wireless Users

Let N be the set of satisfied users whose minimum data requirements, D_{min} , have been met at the end of a schedule cycle. The primary objective is *to maximize the size, $|N|$, in each cycle*. In addition, the secondary objective is to maximize the utility of those $|N|$ users and hence the system throughput. After the scheduler has allocated time slots to the users, there might exist *residual slots* which are insufficient to satisfy the minimum data requirement of any additional unallocated user. In order to maximize the system throughput, these slots are distributed among the allocated users. Thus, the *Optimal Scheduling Policy* can be constructed as follows:

$$\text{maximize } |N| \tag{4.1}$$

$$\text{such that } \begin{cases} \sum_{j=1}^n w_{i,j} = m \\ U_j \geq U_{min} \text{ for } 1 \leq j \leq n \\ w_{i,j} \geq 0 \end{cases} \tag{4.2}$$

In order to achieve the optimal schedule, we first formulate the problem in terms of linear programming (LP). In the next section, we show that the LP is equivalent to the optimal scheduling policy.

4.1.6 LP Formulation

Considering that $w_{i,j}$ determines the schedule matrix and $r_{i,j}$ the data rate of user j for slot i , the system throughput (T_p) for a schedule cycle is given by:

$$T_p = \sum_i \sum_j w_{i,j} r_{i,j} \quad (4.3)$$

However, since the optimal schedule does not maximize the throughput for each slot, the system suffers throughput loss governed by a penalty function Θ_P . The penalty function measuring the system throughput loss for every schedule cycle is given by:

$$\Theta_P = \sum_i \sum_j \left((\max_j r_{i,j}) - w_{i,j} r_{i,j} \right) \quad (4.4)$$

Consequently, the total utility Θ_U , derived by the users is given by:

$$\Theta_U = \sum_{j=1}^N U_j \quad (4.5)$$

where U_j denotes the utility of user j as a function of the data received. Since the effective objective function is to obtain the joint performance measure of all the user utilities as well as the system throughput loss, we employ a value function, V , to map *both* the penalty and the user utility to a common unit (e.g., money metric) so that the joint optimization can be achieved. We define the value penalty (V^P) as a function of Θ_p and the value user satisfaction (V^U) as a function of Θ_U . Thus,

$$V^P = f(\Theta_P) \quad (4.6)$$

$$V^U = g(\Theta_U) \quad (4.7)$$

For the time being, let us ignore the specific nature of $f(\cdot)$ and $g(\cdot)$. Consequently, the overall objective function of the system, the optimization of which would provide the solution to expression (4.1), can be written as:

$$\textit{System Objective: Maximize } (V^U - V^P) \quad (4.8)$$

subject to the conditions stated in expression (4.2). Note that in the process we have mapped inequality (4.1) to an alternative formulation given by expression (4.8). Next we model the user utility functions.

4.2 Modeling User Utility

The methodology of quantifying the user satisfaction derived from the services received using the concept of *utility functions* has been established in [45]. We assume that the utility is an increasing function of the data received (D_r). However, the utility remains zero unless and until a minimum amount of data (D_{min}) is received, i.e., even for non-zero D_r , the utility is zero if $D_r < D_{min}$. This can be justified by the fact that most applications require a minimum amount of data below which the applications fail to execute. For example, for streaming multimedia, the media player needs to wait for a certain number of packets before the media frame can be successfully constructed. Formally, we define the utility function for user j receiving D_r amount of data as follows:

$$U_j(D_r) = \begin{cases} 0 & 0 < D_r < D_{min} \\ U_j(D_r) & D_{min} \leq D_r < D_{max} \\ U_{max} & D_r \geq D_{max} \end{cases} \quad (4.9)$$

We consider a generic utility function as shown in Figure 4.2. Clearly the change in utility is more prominent between D_{min} and D_{max} . This kind of utility function is very intuitive and can be better illustrated by the following example. Consider that a video demands somewhere between 1 Mbps and 4 Mbps of bandwidth. This means that with an effective bandwidth of less than 1 Mbps, the quality of the video is too poor to be perceivable. On the other hand, with an effective bandwidth of more than 4 Mbps, there is no perceptible improvement in the video quality. Thus, in Figure 4.2, $D_{min} = 1$ Mbps and $D_{max} = 4$ Mbps. The corresponding utilities derived by the user are U_{min} and U_{max} , respectively. To satisfy U_{min} , the corresponding resource (i.e., the number of slots) must be available. Moreover, a smaller U_{min} does not necessarily mean that the number of slots required will be less since the instantaneous channel conditions might be bad; thus requiring more slots to provide the minimum utility. Beyond D_{min} , we consider the user utility function along the lines of diminishing returns, i.e., the *marginal utility*[†] of the user diminishes as a function of the allocated data. Additionally, the marginal utility is zero or negligible when $D_r > D_{max}$. In other words, the utility does not significantly increase if received data exceeds D_{max} , as illustrated in Figure 4.2. Also, in all our analysis the user utility has been normalized between 0 and 1 where U_{min} and U_{max} are the corresponding threshold utilities.

4.2.1 Utility Function and the Exposure Problem

Let us now understand the inter-dependencies between the utility function and the exposure problem. According to the auction theory terminology [68], the exposure problem arises because the users' valuations for the number of available slots are not additive. This implies that there exists *complementarity* or *substitutability* among the slots.

[†]In economics, "marginal utility" is the additional utility (satisfaction or benefit) that user derives from an additional unit of service such as time slot, in our case.

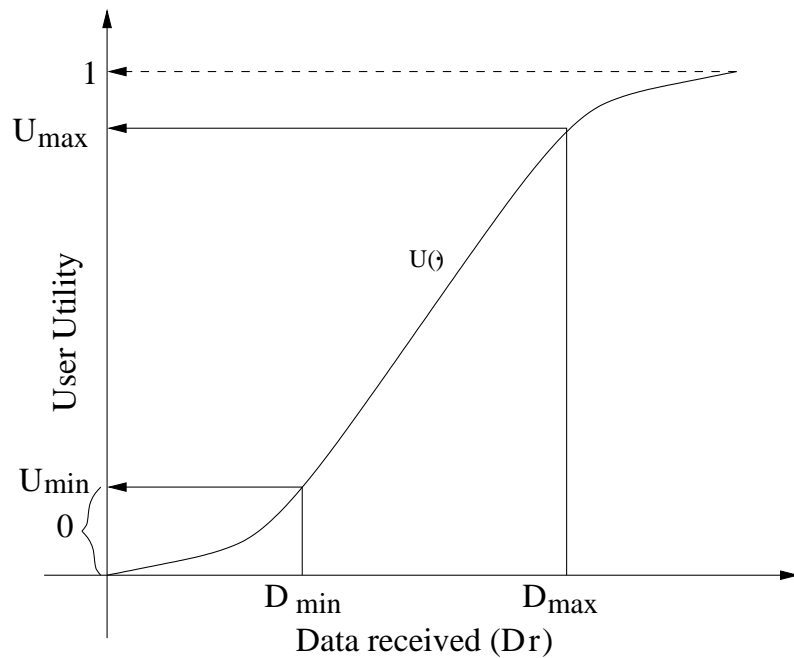


Figure 4.2. Utility curve illustrating the marginal utility of the user as a function of the data received.

Although a user might be allocated slots according to the wireless channel condition as in opportunistic scheduling, or some fixed number of slots based on temporal fairness, the minimum data requirements might not be satisfied. In auction terminology, the user might be enticed to bid a higher price for a subset of the desired bundle with the hope to acquire the total bundle, but ends up gaining nothing since the minimum requirement is not satisfied.

We first define and then employ the complementarity and substitutability effects to demonstrate that the exposure problem depends on whether the user utility function is linear or non-linear. As outlined in Section 4.1.3, the exposure problem helps us to identify if the slot allocation (i.e., the scheduling) is being performed effectively. In the complementarity effect, the value maximization for the system is only achieved by allocating a particular bundle of slots but not any subset of it. The substitutability effect encompasses the scenario when the value maximization of the system is achieved only

when the right bundle of slots is allocated and not any superset of it. Note that, opportunistic scheduling does not solve the exposure problem for non-linear utility functions. This is because a non-linear utility function, κ , displays sub-additivity or super-additivity over various ranges, i.e.,

$$\kappa(x) + \kappa(y) \leq \kappa(x + y) \quad \text{or} \quad \kappa(x) + \kappa(y) \geq \kappa(x + y) \quad (4.10)$$

Opportunistic scheduling mechanisms concentrate on the *current slot* to be scheduled and base their decision on an objective function. These schemes do not consider prior or future allocations, and thus are unable to capture the complementarity and substitutability effect between the slots.

4.3 Scheduling and Combinatorial Auctions

In this section, we highlight the equivalence between the optimal resource allocation problem in multi-rate wireless systems and combinatorial auctions and by deriving the mapping between the two. When the objective in a market is achieved, such as value maximization (resp. minimization) for seller (resp. buyer) in forward (resp. reverse) auctions, the market is said to be in an *equilibrium* state. The market equilibrium corresponds to the optimal schedule as defined in Section 4.1.5 where the goods map to the time slots and the objective is to satisfy expression (4.8). Under such circumstances, the combinatorial auction problem can be formulated as follows.

Let $M = \{1, 2, \dots, m\}$ denote the set of goods available for auction and let $u_j(S)$ denote the utility a seller derives if the buyer j acquires the bundle S . Consequently, the utility is formulated as:

$$u_j(S) = \sum_{x \in X} \beta_x V_{j,x}(S_b) \quad (4.11)$$

where X is the set of factors determining the overall utility of a bundle S_b ; β_x is the weight for a given factor x ; and $\sum_{x \in X} \beta_x = 1$. The term $V_{j,x}(S)$ is the value of the factor x by allocating the bundle S_b to buyer j . We define $\psi(S_b, j)$ as:

$$\psi(S_b, j) = \begin{cases} 1 & \text{if bundle } S_b \text{ is allocated to buyer } j \\ 0 & \text{otherwise.} \end{cases} \quad (4.12)$$

Thus, the forward combinatorial auction can be formulated as the following optimization problem:

$$\text{maximize } \sum_{j \in N} \sum_{S_b \subseteq M} u_j(S) \psi(S_b, j) \quad (4.13)$$

$$\text{such that } \begin{cases} \sum_{i \ni S_b} \sum_{j \in N} \psi(S_b, j) \leq 1 & \forall i \in M \\ \sum_{S_b \subseteq M} \psi(S_b, j) \leq 1 & \forall j \in N \\ \psi(S_b, j) = \{0, 1\} & \forall S_b \subseteq M, \forall j \in N \end{cases} \quad (4.14)$$

The first condition ensures that the overlapping sets of items are never assigned, whereas the second one ensures that no bidder receives more than one subset. The reverse combinatorial auction can be formulated in a similar fashion. Note that in the reverse auction there exists a single buyer intending to procure items from multiple sellers who quote the minimum price. In both the forward and reverse scenarios, it is assumed that the slots display complementarity and substitutability in terms of utility and costs, respectively. A careful observation of expression (4.13) reveals that the formulation is identical to the system objective problem defined in expression (4.8). Finding the solution to expres-

sion (4.13) is known as the *winner determination problem* [68] for both combinatorial forward and reverse auctions.

4.3.1 Forward Auction for Multi-Rate Slot Allocation

Consider the scenario where the wireless system represents the seller and the users the buyers. Recall that the primary system objective of scheduling is to maximize the number of users whose minimum utility, D_{min} , is satisfied. The secondary objective is to maximize the system throughput once it is no longer possible to add a user whose D_{min} can be satisfied. Such conditions require that the auction be done in *two stages*. In the first stage, bids satisfying the minimum utility are determined. Let \mathcal{B}_I be the initial feasible bid set for the forward auction. It is defined as the set of all sets (bids) of slots which are maximal sets satisfying the minimum utility for all the users. The *initial phase* is thus described by:

$$\text{maximize } \sum_{j \in N} \sum_{S_b \in \mathcal{B}_I} u_j(S_b) \psi(S_b, j) \quad (4.15)$$

$$\text{such that } \begin{cases} \sum_{i \in S} \sum_{j \in N} \psi(S_b, j) \leq 1 & \forall i \in \mathcal{B}_I \\ \sum_{S_b \in \mathcal{B}_I} w(S_b, j) \leq 1 & \forall j \in N \\ \psi(S_b, j) = \{0, 1\} & \forall S_b \in \mathcal{B}_I, \forall j \in N \end{cases} \quad (4.16)$$

The *termination criterion* for the first phase of slot allocation occurs when no additional user can be granted D_{min} amount of data. The *second phase* consists of allocating the residual slots which are available at the end of first phase of the scheduling operation. Depending on the objective and utility functions, either a second round of auctions or any standard opportunistic scheduling algorithm can be employed to disseminate the

residual slots among the users. Though there exists a solution for the above formulation, finding the set of winners is shown to be an \mathcal{NP} -complete problem [32] and cannot even be approximated to a ratio of $n^{1-\epsilon}$ in polynomial time where n is the total number of users and $0 < \epsilon < 1$. As a result, it is infeasible to implement forward combinatorial auction based scheduling for real-time multi-rate wireless systems. This motivates us to explore reverse auction based scheduling in the next section.

4.4 Multiple Slot Scheduling through Reverse Auctions

In this section, the delay-sensitive multi-rate scheduling problem is reformulated based on *reverse combinatorial* auction. In such a scenario, the wireless base station is the buyer who wants to procure m slots and the set of N users are the sellers each having m slots of different values (a.k.a. data rates). The prices that the users quote for the bundle of slots depend on the utility derived by the user when the base station procures those slots. The problem can be formally stated as:

$$\text{minimize } \sum_{j \in N} \sum_{S_b \subseteq M} p_j(S_b) \psi(S_b, j) \quad (4.17)$$

$$\text{such that } \begin{cases} \sum_{i \in S_b} \sum_{j \in N} \psi(S_b, j) \geq 1 & \forall i \in M \\ \sum_{S_b \subseteq M} \psi(S_b, j) \leq 1 & \forall j \in N \\ \psi(S_b, j) = \{0, 1\} & \forall S_b \subseteq M, \forall j \in N \end{cases} \quad (4.18)$$

Here $p_j(S_b)$ denotes the price that user j quotes for the slot bundle S_b . The solution to the above problem is nothing but the *winner determination problem*. In this framework, the users compete against each other to sell the set of slots to the base station. They are

deprived of some value if they cannot get the base station to buy the slots from them. Identifying the *deprivation function* is essential for deciding the best set of slots for any user. Thus, the quoted price for any bundle is a function of the deprivation function.

4.4.1 Deprivation Function

The objective of the user is to dispose the set of slots and obtain the desired utility. The buyer (base station) buys the set of slots only if a specific minimum value of D_{min} is achieved during the bidding process. Failure to sell the slots deprives the user of the minimum utility. Hence, a deprivation value is associated with the set of slots not acquired by the base station from that user. The deprivation function depends on two factors :

- (i) The utility derived by the user by giving the bundle to the base station.
- (ii) The throughput loss the base station may experience while procuring the bundle.

The utility that the user gets from a bundle of slots, S_b , depends on the type of applications. Following expression (4.9), the utility function can be defined as: $V_j^U = V(U_j)$, where $V(\cdot)$ is a value function mapping the utility to equivalent money metric. Similarly, the monetary equivalent of throughput loss that the system experiences by acquiring bundle S_b from user j is calculated using expression (4.6) and is denoted by $V_j^L(S_b)$. Therefore, we can define the deprivation function for a slot bundle S_b as:

$$P_j(S_b) = \alpha V_j^U(S_b) + (1 - \alpha) V_j^L(S_b) \quad (4.19)$$

Here β is a control or tunable design parameter that controls the relative weight of the two attributes. For $\beta = 1$, the deprivation function basically boils down to guaranteeing only user satisfaction, whereas for $\beta = 0$, the system considers only throughput maximization.

4.4.2 Mechanism for Reverse Auction

We use the simple *single-round, sealed-bid first* price combinatorial reverse auction mechanism. All the “asks” (or quotes) are submitted prior to a deadline and the slot allocation is achieved based on the set of “asks” received. The throughput would have been drastically penalized had the auction been non-incentive based. However, the equilibrium is *not* guaranteed for non-incentive combinatorial auctions [64]. In general, whether the mechanism is incentive compatible or not, the price $p_j(S)$ quoted by user j for the bundle S_b is a function of the deprivation function P_j . Thus:

$$p_j(S_b) = f_j(S_b)P_j(S_b) \quad (4.20)$$

where $f_j(S_b)$ is the *price mapping function* that defines the relationship between the price and the deprivation function. Since we have assumed incentive compatible auction mechanism, $f_j(S_b) = -1$, for all j and for all S_b .

The solution to the winner determination provides the desired schedule vector. It has been shown in [32] that in reverse auction, approximate solutions can be developed inspite of the fact that the problem is \mathcal{NP} -complete. Hence, we develop our slot procurement algorithm along the lines of reverse auction. However, since the primary and secondary objectives have conflicting goals, we decouple the algorithm into two phases. In the first phase, the *restricted* phase, we compute the set of users whose minimum requirement is satisfied. That is, slots are acquired from as many users as possible while requiring that each user is able to get rid of the minimum deprivation value. During the second phase, the *unrestricted* phase, the residual slots are allocated which cannot satisfy D_{min} for any additional user. These two phases are described below.

4.4.3 Restricted Phase

In this phase, multiple single-round reverse auctions are held till no additional user is able to get rid of the minimum deprivation value. In each round, the system considers “asks” from the users on the remaining minimal unallocated bundles. This means that the bundle should just be able to get rid of the minimum deprivation value. Each user is allowed to provide an “ask” for only one bundle of slots A_j . From these initial “asks”, the initial feasible bundle or set (\mathcal{A}^k) is constructed for round k of the restricted phase.

For each round, the reverse auction takes the following formulation:

$$\text{minimize } \sum_{j \in N} \sum_{S_b \in \mathcal{A}^k} p_j(S_b) \psi(S_b, j) \quad (4.21)$$

$$\text{such that } \begin{cases} \sum_{i \ni S_b} \sum_{j \in N} \psi(S_b, j) \geq 1 & \forall i \in \mathcal{A}^k \\ \sum_{S_b \in \mathcal{A}^k} \psi(S_b, j) \leq 1 & \forall j \in N \\ \psi(S_b, j) = \{0, 1\} & \forall S_b \in \mathcal{A}^k, \forall j \in N \end{cases}$$

Once the minimum deprivation value of a user has been satisfied in a certain round, the user is barred from taking part in subsequent rounds of the auction process in the restricted phase. Let the ACQUIREDSET denote the set of accepted “asks” and each “ask” A_j is represented by a set of vector $\langle \lambda_1^j, \lambda_2^j, \dots, \lambda_m^j \rangle$, where λ_i^j is 1 if the i th slot is in the “ask” for user j , otherwise it is 0. Let PERMITTEDSET be the set of permitted “asks”. Let us define θ_i such that $\theta_i = 1$ if the i th slot has not been acquired, otherwise $\theta_i = 0$. Let SATISFIEDSET be the set of users whose minimum deprivation value has been satisfied. The algorithm for the restricted phase is described in Figure 7.

Algorithm 7 Restricted Phase: Algorithm for Winner Set Determination

- 1: Initialize ACQUIREDSET = SATISFIEDSET = ϕ , the current round $k = 0$ and for
all $i \in M$, $\theta_i = 1$
 - 2: Construct \mathcal{A}^k from $j \in (N - SATISFIEDSET)$
 - 3: **while** $\mathcal{A}^k \neq \phi$ **do**
 - 4: Initialize PERMITTEDSET = \mathcal{A}^k
 - 5: **while** PERMITTEDSET $\neq \phi$ **do**
 - 6: Find j^* such that $p_j(S)$ is minimum and $A_j \in PERMITTEDSET$
 - 7: **if** A_{j^*} has a slot i , for which $\lambda_i^{j^*} = 1$, but $\theta_i = 0$ **then**
 - 8: Remove A_{j^*} from the PERMITTEDSET
 - 9: **else**
 - 10: Add A_{j^*} to ACQUIREDSET
 - 11: Remove A_{j^*} from the PERMITTEDSET
 - 12: For all $i \in A_{j^*}$, make $\theta_i = 0$
 - 13: Add j^* to SATISFIEDSET
 - 14: **end if**
 - 15: Increment k by 1
 - 16: Construct \mathcal{A}^k from $j \in (N - SATISFIEDSET)$
 - 17: **end while**
 - 18: **end while**
 - 19: Return ACQUIREDSET
-

4.4.4 Unrestricted Phase

The residual slots aid in achieving the secondary objective of maximizing the utility of allocated users as well as maximizing the system throughput during the unrestricted

phase. The allocated users strive to further minimize the deprivation value by selling their slots. However, unlike the restricted phase, there is no restriction on the size of slot bundle. Note that none of the users whose minimum utility (i.e., the minimum deprivation value) has not been satisfied, is allowed to compete in this phase. Additionally, the slots may not exhibit complementary/substituability relationship. Consequently, the exposure problem explained earlier will not occur. Under such conditions when the utility is assumed to be linear, scheduling the residual slots can be performed by employing any one of the existing opportunistic scheduling algorithms.

On the contrary, if complementary/substituability effect exists between the residual slots, the allocation should be performed using combinatorial reverse auction so as to overcome the exposure problem. For the unrestricted phase, the "asks" are based on the further reduction of the deprivation value. The objective for the buyer (i.e., the system), is now set to choose the "asks" from the users, which minimizes its total price. This guarantees throughput maximization for both the system as well as the chosen users. The auction proceeds similar to the restricted phase but continues till all the slots have been exhausted. Let \mathcal{A}^k denote the set of all remaining slots after round k of the unrestricted phase. The mathematical formulation for round k is given by:

$$\text{minimize } \sum_{j \in N} \sum_{S_b \in \mathcal{A}^k} p_j(S_b) \psi(S_b, j) \quad (4.22)$$

$$\text{such that } \begin{cases} \sum_{i \in S_b} \sum_{j \in N} & \psi(S_b, j) \geq 1 \quad \forall i \in \mathcal{A}^k \\ \sum_{S_b \in \mathcal{A}^k} & \psi(S_b, j) \leq 1 \quad \forall j \in N \\ \psi(S_b, j) = \{0, 1\} & \forall S_b \in \mathcal{A}^k, \forall j \in N \end{cases} \quad (4.23)$$

Note that the difference between expressions (4.22) and (4.23) is the type of “asks” possible and the set of slots which are part of the reverse auction. The ACQUIREDSET obtained in the previous algorithm is used in the Unrestricted phase. The algorithm is described in Figure 8.

Algorithm 8 Unrestricted Phase: Residual Slot Allocation Algorithm

- 1: Initialize $l = 0$, PERMITTEDSET = ϕ
 - 2: Construct \mathcal{A}^l from $j' \in \text{SATISFIEDSET}$ and slots for which $\theta_i = 1$
 - 3: **while** $\theta_i = 1$, for some $i \in M$ **do**
 - 4: Initialize PERMITTEDSET = \mathcal{A}^l
 - 5: **while** PERMITTEDSET $\neq \phi$ **do**
 - 6: Find j'^* such that $A'_{j'} \in \text{PERMITTEDSET}$ and $p_j(S)$ is minimum
 - 7: **if** $A'_{j'^*}$ has a slot i , for which $\lambda_i^{j'^*} = 1$, but $\theta_i = 0$ **then**
 - 8: Remove $A'_{j'^*}$ from the PERMITTEDSET
 - 9: **else**
 - 10: Add $A'_{j'^*}$ to ACQUIREDSET
 - 11: Remove $A'_{j'^*}$ from the PERMITTEDSET
 - 12: For all $i \in A'_{j'^*}$, make $\theta_i = 0$
 - 13: **end if**
 - 14: Increment l by 1
 - 15: Construct \mathcal{A}^l from $j' \in \text{SATISFIEDSET}$ and $\theta_i=1$
 - 16: **end while**
 - 17: **end while**
 - 18: Return ACQUIREDSET
-

After the execution of the unrestricted phase algorithm, the ACQUIREDSET is updated which provides the distribution of the slots for the schedule cycle under consideration.

4.5 Performance Analysis

In this section, we analyze the proposed algorithms.

Theorem 1. *The worst case running time for the restricted phase slot procurement algorithm is $\mathcal{O}(n^2m)$ where m is the number of slots in each schedule cycle and n is the number of users.*

Proof. Assume $m \gg n$. The complexity of the algorithm in the restricted phase depends on the "ask" construction (line 2) in Figure 3 and the selection of appropriate j (line 6). In the worst case, line 2 of the algorithm takes $(n - k)(m - k)$ operations where k is the current round. Line 6 takes $n - k$ operations in the worst case. The maximum number of possible rounds is n . So, the worst case complexity can be given by $\mathcal{O}(\sum_{k=0}^{n-1} [(n - k)(m - k) + (n - k)]) = \mathcal{O}(\sum_{k=0}^{n-1} [nm - (m + 1)k + k^2]) = \mathcal{O}(mn(n - 1) + (m + 1) \sum_{k=0}^{n-1} k + \sum_{k=0}^{n-1} k^2) = \mathcal{O}(n^2m) + \mathcal{O}(n^2m) + \mathcal{O}(n^3) = \mathcal{O}(n^2m)$. since $n \ll m$; \square

Corollary 1. *The worst case complexity of the unrestricted phase slot procurement algorithm is $\mathcal{O}(n^2m)$.*

Lemma 1. *Let the effective price for each slot be defined as $\mathcal{P}(o) = \frac{p_j}{|A_j|}$ where p_j is the price paid by user j for acquiring the set of slot A_j . If OPT is the total cost that the base-station pays for the optimal solution, then*

$$\mathcal{P}(o_k) \leq \frac{OPT}{l' - k + 1}$$

where $\{o_i\}, i = 1, \dots, l$ is an ordering of the slots based on the sequence in which they are acquired by the base station, l is the total number of slots procured by the base station in the restricted phase and $l' \leq l$.

Proof. Let o_k be covered (i.e., these slots are taken up) when the “ask” A_j was picked by the algorithm. After, o_k , there are at least $l' - k + 1$ slots to be covered. Since the optimal cost OPT covers all the l slots, it can also cover the remaining $l' - k + 1$ slots. So, there must be at least one “ask” whose average cost of covering is at most $\frac{OPT}{l' - k + 1}$. As our algorithm chooses the slots from the lowest to the highest average cost per slot, $\mathcal{P}(o_k) \leq \frac{OPT}{l' - k + 1}$. \square

Theorem 2. *The restricted phase slot procurement algorithm finds a solution that is within a factor $(1 + \log m)$ of the optimal solution where m is the total number of slots to be procured.*

Proof. The proof for the bound is similar to the one presented in [68]. Let l be the total number of slots that could be covered by the optimal solution and l' be the total number of slots that are covered by our algorithm. From Lemma 1, the proof of Theorem 2 can be outlined as follows: Let the “asks”, that were picked in the restricted phase that are able to get rid of the minimum deprivation value be denoted by $A_{j_1}, A_{j_2}, \dots, A_{j_s}$, where j_s is the last user whose bundle of slots was chosen. The total cost is given by $\sum_{x=1}^s p_{j_x} = \sum_{k=1}^{l'} \mathcal{P}(o_k)$. Using Lemma 1, the total cost can be written as:

$$\sum_{k=1}^{l'} \mathcal{P}(o_k) \leq OPT \left(1 + \frac{1}{2} + \dots + \frac{1}{l'}\right) \leq OPT \times H_{l'}$$

where $H_{l'}$ is the l' th harmonic number. Since

$$H_{l'} \leq 1 + \ln l' \leq 1 + \ln l \leq 1 + \ln m,$$

the cost is bounded by $(1 + \log m)$ of the optimal. \square

Lemma 2. *Let \tilde{n}_i be the number of slots obtained by user i in order to satisfy the minimum utility using our scheme and \hat{n}_i be the number of slots obtained using an opportunistic scheme. Then $\hat{n}_i \geq \tilde{n}_i, \forall i \in N$.*

Proof. In our scheme, the slot allocation always tries to give the best available slots to any user so as to satisfy the minimum utility requirement at the minimum cost. Here \tilde{n}_i is the minimum number of slots required to satisfy the minimum utility. Now consider an opportunistic scheme where the decision is based on a slot by slot basis. Consider a user whose minimum utility has been satisfied. If the user's best available slots come in descending order of their individual utility values, then the user will reach the minimum utility level with the smallest number of slots. In this case, $\hat{n}_i = \tilde{n}_i$. Otherwise, the user may get another slot which is not the user's available slot. Therefore, to satisfy the minimum utility, the user will require at least the minimum number of slots. In either case: $\hat{n}_i \geq \tilde{n}_i$. \square

Theorem 3. *Let N_c be the set denoting the maximum number of users whose minimum utility has been satisfied by the combinatorial reverse auction based scheduling and let N_o be the set of users who have been satisfied by the opportunistic scheme. Then $|N_c| \geq |N_o|$.*

Proof. From Lemma 2, $\hat{n}_i \geq \tilde{n}_i$, for all i whose minimum utility have been satisfied. By contradiction, let us assume $|N_c| < |N_o|$. Then

$$\sum_{i=1}^{|N_o|} (\hat{n}_i + \hat{k}_i) + l = \sum_{i=1}^{|N_c|} (\hat{n}_i + \tilde{k}_i) \quad (4.24)$$

where \hat{k}_i and \tilde{k}_i are the extra slots given to user i after satisfying the minimum utility, and l is the total number of slots given in the case of opportunistic scheduling to users whose minimum utility could not be satisfied. The above equation can be rewritten as

$$\sum_{i=1}^{|N_c|} (\hat{n}_i - \tilde{n}_i) + \sum_{i=|N_c|+1}^{|N_o|} \hat{n}_i + l + \sum_{i=|N_c|}^{|N_o|} \hat{k}_i = \sum_{i=1}^{|N_c|} (\tilde{k}_i - \hat{k}_i) \quad (4.25)$$

But this implies that:

$$\sum_{i=1}^{|N_c|} \tilde{k}_i \geq \sum_{i=|N_c|}^{|N_o|} \hat{n}_i \quad (4.26)$$

This is clearly not possible since it would mean that our auction based scheme would be able to accommodate at least one more user using the \tilde{k}_i 's. Hence, $|N_c| \geq |N_o|$. \square

4.6 Simulation Study

This section studies the effectiveness of our proposed scheduling scheme through simulation experiments. We also compare how the auction based scheme fares with respect to two extreme scheduling disciplines: round-robin and throughput maximization - that serve as the basis for comparing the fairness and maximum system throughput, respectively. We study how each scheme performs in terms of the number of satisfied users and global system throughput.

4.6.1 System and Channel Model

We consider a single cell wireless data network for our simulation study due to the fact that the scheduling schemes under evaluation are designed to work best in the presence of a single base station. We also assume that all the users under consideration are receiving real-time streaming multimedia traffic. In order to support multimedia traffic (MPEG-4 or H.263) of various qualities (low, medium, and high) as given in

[28], we consider three values for D_{min} : 16 Kbps, 64 Kbps, and 128 Kbps. We model our simulation based on the HDR system that is capable of supporting 11 different data rates with each schedule cycle consisting of 1000 slots. We assume user mobility is random (both speed and direction) and employ the path-loss model and the slow log-normal model [42] for wireless channels.

4.6.2 Simulation Results

For our proposed auction based scheduling scheme, the variation of system throughput with the number of users for different values of D_{min} is shown in Figure 4.3. As expected, the system throughput initially increases but ultimately gets saturated with the increase in the number of users. Next, we identify the maximum system capacity

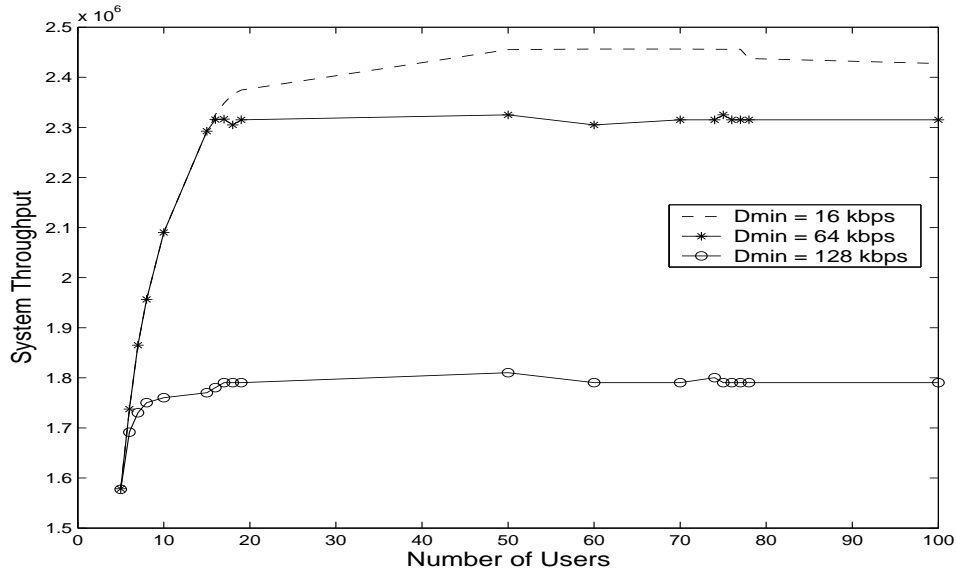


Figure 4.3. Throughput vs. Number of Users in the System. Notice how the throughput decreases with increase in D_{min} .

in terms of satisfied users by setting D_{min} to different values. For each value of D_{min} , we obtained a range of users who are satisfied by the system. This is shown in Table 4.2.

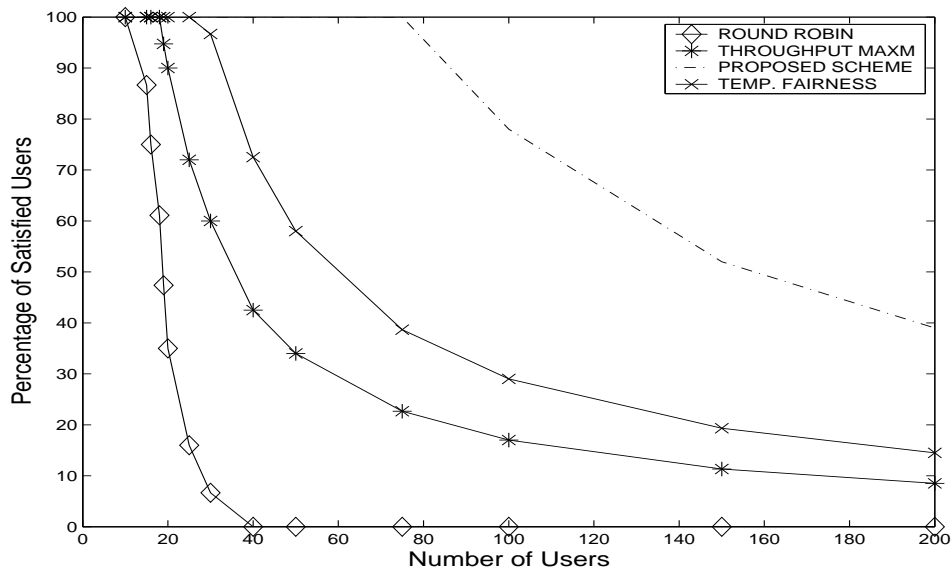


Figure 4.4. Performance of each scheduling scheme measured using satisfied users as a percentage of the total users.

It is logical that for smaller D_{min} , a greater number of users can be satisfied. The

Table 4.2. D_{min} vs. Maximum Number of users

D_{min} (kbps)	Max Satisfied Users
16	76-82
64	16-20
128	5-8
256	1-2
512	0-1

comparison of the system throughput achieved by various schemes is shown in Figure 4.5. As expected, the system throughput is the best for the throughput maximization scheme and worst for the round-robin scheduling algorithm. In the case of opportunistic scheduling with temporal fairness, the throughput is penalized. The throughput performance of the proposed auction based scheme is better than both the round-robin and opportunis-

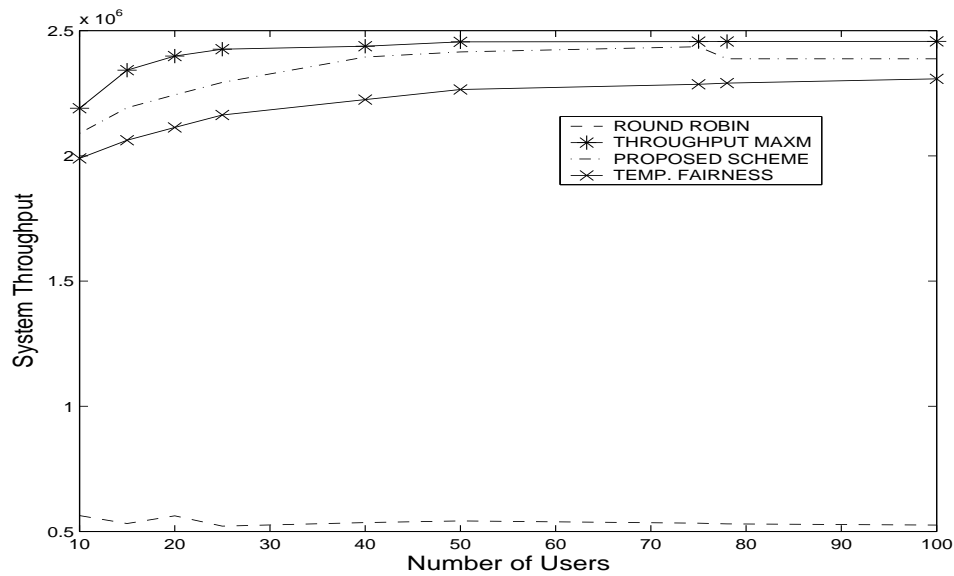


Figure 4.5. Throughput vs. the number of users in the system for different scheduling algorithms.

tic scheduling with temporal fairness, and is very close to the throughput maximization scheme. In order to visualize the working of the proposed scheduling scheme, we consider a hypothetical scenario with 15 users in the system. The users are represented as u_i in Figure 4.6. For the purpose of explaining the workings of our algorithm, a temporal snapshot of four successive scheduling decision cycles is also presented.

With $D_{min} = 128$ Kbps, the first three schedule cycles yield the schedule vector as $[1, 2, 4, 5, 6, 7, 10]$. But for the 4th cycle, user 7 is replaced by user 9. Although the allocated users are receiving D_{min} , the variation of the slot distribution between users is due to the varying channel condition. The scheduling scheme judiciously distributes the residual slots after the restricted phase and does not allocate any more slots if D_{max} is attained, as is the case with user 10 in this example. Careful observation reveals that though all the allocated users were receiving D_{min} or more data, user 7 was receiving lesser slots in each succeeding schedule cycle such that in the 4th cycle, user 7 was eliminated by user 9 in the restricted phase. Thus, the scheduler is intelligent enough to identify

and allocate the user to achieve the system objective. In each schedule cycle, all the allocated users are guaranteed D_{min} amount of data.

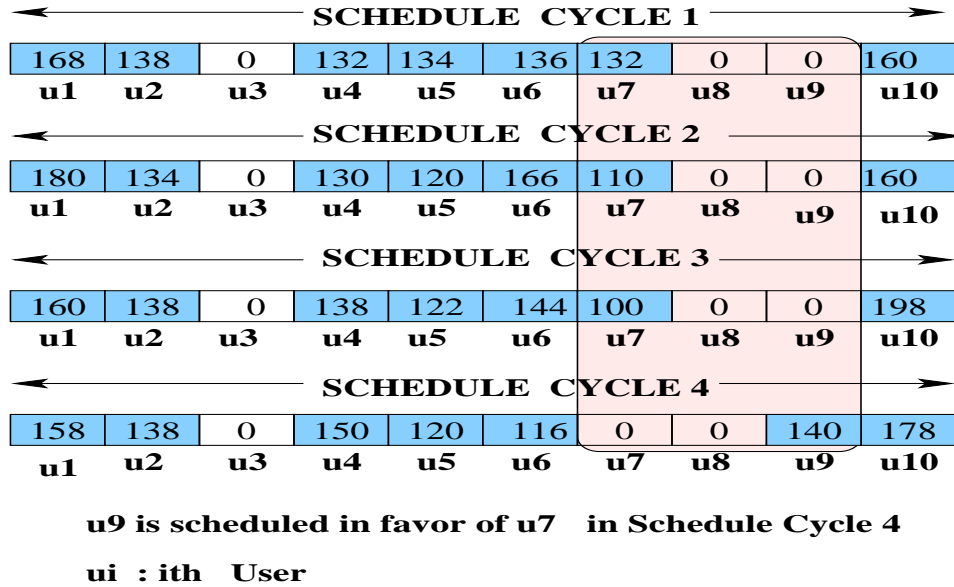


Figure 4.6. Slot Distribution of Users in Schedule Cycle.

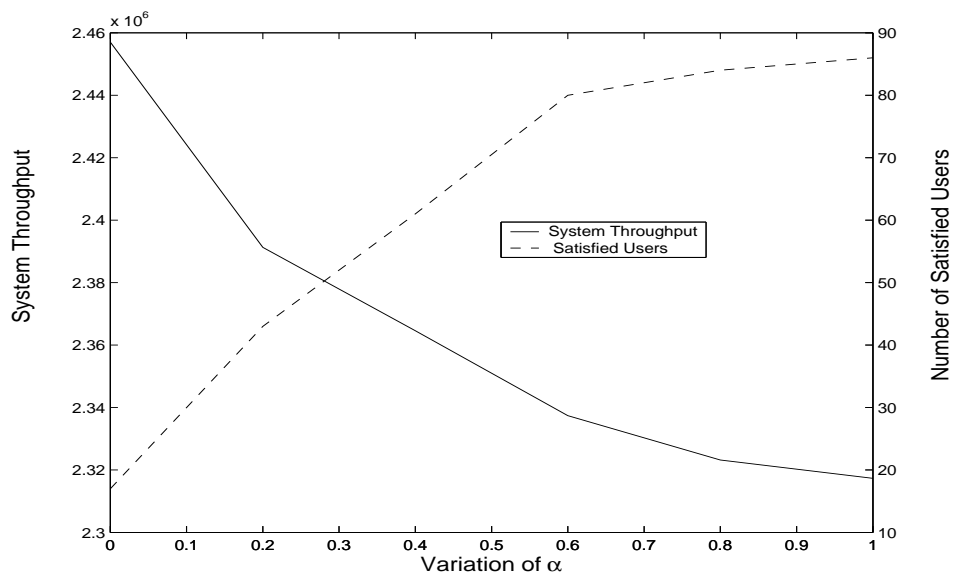


Figure 4.7. System Throughput and Number of Satisfied users vs. α .

Next, we investigate the variation of system throughput and the number of satisfied users with the tunable parameter α , as defined in Equation (4.19). The value of α depends on the objective of the wireless service providers that either maximizes the throughput or guarantees user utility or a combination of both. Hence, we evaluate the system throughput and the number of satisfied users by varying α from 0 to 1. For $\alpha = 0$, the deprivation function becomes totally a function of the throughput maximization whereas for $\alpha = 1$ the deprivation function only cares about the user satisfaction. As expected, the throughput maximizes for $\alpha = 0$, whereas the number of satisfied users is maximized for $\alpha = 1$, as illustrated in Figure 4.7.

4.7 Summary

We have used an auction based scheduling algorithm for allocating the slots in a time-division, multi-rate wireless system and have shown that it is possible to increase the number of satisfied users. We have justified that opportunistic scheduling algorithms that aim to maximize the system throughput are unable to address the exposure problem. We have formalized the slot allocation problem in the form of a market where multiple users bid for the number of slots to satisfy their minimum QoS requirements. With the help of combinatorial auctions, we have shown how the exposure problem can be successfully eliminated. In the process, we have been able to achieve the primary objective of maximizing the number of users whose minimum slot requirements are satisfied. The remaining slots, if any, are allocated with a view of maximizing the system throughput. In our study, we have applied reverse auction theory in order to deal with the real-time scheduling requirements. We have derived the approximation ratio of the auction based scheduling algorithm which results in more satisfied users than other existing opportunistic scheduling schemes. In summary, not only we have proposed resource management algorithms and channel estimation techniques and ensured the user satisfaction but also

we have performed an analysis of how many satisfied users can be supported. However, user satisfaction is closely connected to the network connectivity issues. In the last part of the dissertation, we focus on solutions which ensure network connectivity in multi-rate systems , specifically for wireless LANs.

CHAPTER 5

EMANCIPATING THE IEEE 802.11 NETWORK FROM HANDOFF DELAY

In existing cellular networks infrastructure, mobility management solutions are well established. Due to the presence of control channels, the mechanisms for handling seamless connectivity are governed completely by the base stations. However, in the case of IEEE 802.11 systems, mobility solutions are still a major challenge as the burden of initiating and maintaining seamless connection lies completely with the end hosts. In the previous chapters, we focussed on channel state estimation, smart scheduling as well as rate control algorithms to boost throughput and performance of the system. The algorithms presented so far are applicable to any centralized multi-rate wireless systems. However, mobility solutions in IEEE 802.11 are distributed as mobility management is initiated and controlled by the end hosts. Considering the fact that in new generation of converged networks mobile devices are expected to be paired with both cellular and IEEE 802.11 systems, in this chapter we describe the design, evaluation and implementation of a new framework that facilitates *seamless* and *transparent handoff* between different access points (APs) specifically in standard IEEE 802.11 based wireless local area networks.

Our proposed algorithm is a *client side* solution capable of reducing handoff delays in Wi-Fi networks to 15ms at best, 20ms in the average case and 26ms in the worst case. It is fully compliant with the current IEEE 802.11 standard and has been proven to work with standard IEEE 802.11 compliant wireless network interface cards (WNIC). Unlike existing Wi-Fi handoff solutions, our solution only requires a WNIC with a *single radio* and does *not* require any support from the network infrastructure. We have successfully implemented and tested our proposed algorithm on Atheros AR5212 chipsets using the MadWifi driver and evaluated its performance using different traffic classes in both our educational network at UTA and in controlled wireless testbed.

The remainder of the chapter is organized as follows. In Section 5.1, we motivate the need for fast handoff. We provide experimental results from live networks at UTA in Section 5.2 to elucidate that the beacon inter-arrival process is stochastic. We also expose the dynamics of the underlying process that makes the beacon generation process to be stochastic nature. We believe that this would make deterministic handoff solutions challenging if not impossible in real networks. The design, full implementation insights and working details of the proposed handoff algorithm inside the MadWifi driver [48] running inside the Linux kernel on AR5215 chipset is explained in Section 5.3. We present and analyze the experimental results collected from our wireless testbed at UTA using our framework in Section 5.4. The existing approaches to fast handoff are discussed and compared to in Section 5.5.

5.1 Motivations for Fast Handoff

The huge popularity of Wi-Fi networks is reflected in 2,00,000 Wi-Fi hotspots being deployed worldwide and approximately 200 million 802.11 (Wi-Fi) enabled devices being shipped in 2006 [30]. The success of these 802.11 networks in the infrastructure mode will be highly dependent on the *untethered* network connection provided to the users for the different services offered especially real time services. The bottleneck lies in the fact that typically 802.11 Access Points (AP) have limited range, typically less than 100 meters. A Wireless device, Mobile Node (MN), and Station (STA) will be interchangeably used to refer to a IEEE 802.11 capable device. Since users using such devices are usually mobile in nature, such devices need to quickly attach themselves to another AP as the signal strength from the current AP starts deteriorating. Considering the fact that most real time sessions such as voice, streaming multimedia exchange data in milliseconds time scale, it is important that available handoff algorithms are able to ensure session and network connectivity at such time scales. Hence, the need for efficient

roaming/handoff algorithms.

The inter-packet delay for real time applications typically range from 50ms (VoIP) to around 150ms (less demanding audio/video real-time streaming with frame coding at 64Mbps). Network connectivity is *lost* (VoIP call dropped) or the session *severely degraded* if the STA fails to associate with a new AP within the above mentioned time constraint as it moves out from the vicinity of the currently associated AP. Since application layer performance is what ultimately matters, any successful handoff algorithm must be able to honor such tight timing constraints. Thus the onus lies on the handoff solution to ensure that the STA is seamlessly switched to an AP prior to service degradation. This motivates us for developing a seamless handoff solution.

The IEEE 802.11 standard defined the specifications for 802.11 Medium Access Control (MAC) protocol and RF-oriented PHY parameters, however it has *not* defined any specific roaming (handoff) algorithm and is open to the device vendors to *improvise*. However, at the time of writing, IEEE is working on a new draft called IEEE 802.11r that intends to facilitate Fast Roaming/Transition. The proposal is yet to be ratified by the IEEE committee and it eventually plans to reduce the handoff delays by (i) advocating multi-vendor *compatibility* of wireless hardware devices, and (ii) by defining improved inter-AP communication messages. The 802.11r though minimizes disruption to the MN's datastream but does not allow the client to determine anything about its ability to communicate with other APs over the air. *It is thus not a client end solution.* Hence, the already deployed 802.11 hardware (currently there are approximately 70 million already in use [30]) needs to be upgraded for taking advantage of the newly defined approach. In addition, the decision of (i) when to initiate handoff, (ii) which AP to move to, and (iii) how to manage the handoff process still lies with the MN. It is, however, worth mentioning at this stage that the Inter-Access Point Protocol (IAPP) or IEEE 802.11f is already standardized and it ensures transferring the client security session between

different APs during a handoff process. The IAPP serves to complement our proposed solution by reducing the network authentication time.

5.1.1 Contributions

Having substantiated the need for fast handoff schemes, we delve into the desirable characteristics of such a scheme. In view of the new and legacy IEEE standard, any handoff solution to be successful and widely used must satisfy the following important criteria: (1) should be operable within the IEEE 802.11 standard; (2) must require minimal preferably zero support from the network infrastructure for it to be widely accepted; (3) client-end software solution is desirable so as not to touch the already deployed APs, and hotspots and (4) should be able to determine when to trigger the handoff and most importantly to which AP it should switch to.

We advocate controlled proactive discovery, monitoring and maintenance of neighboring APs, and subsequent handoff execution based on the status of the currently associated AP. Currently, our proposed handoff solution is available as a loadable kernel module for the Linux kernel (2.4.x and 2.6.x) and has been tested to work with dual band WNIC from Netgear (model WAG511) using Atheros Communication AR5212 chipset [7] with absolutely no support from the network infrastructure. The solution achieves 20 ms handoff delays on the average making the MN oblivious to handoff delays. Moreover, in most of the experiments conducted for Layer 2 handoff, we have experienced less than 1% packet loss using commercial off-the-shelf APs. Network connectivity is *lost* (VoIP call dropped) or the session *severely degraded* if the MN fails to associate with a new AP within the above mentioned time constraint as it moves out from the vicinity of the currently associated AP. Since application layer performance is what ultimately matters, any successful handoff algorithm must be able to honor such tight timing constraints. The novel contributions of this chapter are summarized below :

- The design, implementation and evaluation of a handoff framework for Wi-Fi networks (IEEE 802.11a/b/g) which achieves a best case delay of 15ms and guarantees that the delay is bounded between 20ms in the average case and 26ms in the worst case.
- The solution introduces *dynamic* and *adaptive* discovery
- The solution judiciously determines *when* to trigger the handoff to ensure seamless network connectivity as well as to prevent the well known classical ping pong effect.
- Real time client end packet buffering has been implemented to avert packet loss and also to convey the impression to the application layer and to the network stack, the network interface appears to be “on” (even during handoff).

5.1.2 802.11 Handoff: Formal Definition

The process of transferring active session(s) of an MN attached to an AP in a specific radio frequency within a Wi-Fi cell to another AP in the same or different frequency residing in the same or different Wi-Fi cell is defined as *handoff*.

Handoff Phases: A Layer 2 (i.e., no change in IP address) handoff process consists of the following four phases: (i) *scanning* of new APs, (ii) *authentication* with a selected, (iii) *association* with the selected AP, and (iv) *wired update* for packets to be delivered through the new AP. Of the four phases, *scanning* is the most *delay prone* entity with the process typically incurring a delay anything between 350 – 1200ms [70, 71, 65]. However, once an AP has been selected, authentication, association and wired update takes another 30ms [65] before the MN can start communicating with the network. Thus, on an average, the handoff process can take several hundreds of milliseconds to complete.

Different from Cellular Network: At this point the reader might be wondering as to why do not we use approaches used in cellular networks for handling handoffs in Wi-Fi domain? After all both are wireless packet data networks. The answer lies in the dif-

ference in the architecture of Wi-Fi and cellular networks. While in cellular networks the process of handoff is managed by the *network* itself (via specially allocated control channels), the 802.11 standard requires that the *handoff process be managed by the client* without any apriori knowledge of the topology of the network. Further, there are no special control channels available for exchanging network related management information. All these makes the handoff process difficult and different from the traditional cellular domain.

Wi-Fi Handoff - Quantitative Evaluation: The scanning phase poses the biggest hurdle for fast seamless handoff process. In existing approaches during the scanning phase, the MN switches to a particular channel and waits for beacons for AP discovery. The process is repeated for all the channels. The scanning delay can be reduced if *before* the handoff is initiated, the list of accessible APs is already known to the MN. This process, commonly referred to as *background scanning*, can be achieved in two ways: *passively*, when the MN quietly waits for the beacon frames from APs; or *actively*, when *probe request* packets are specifically broadcast and the *probe response* packets from the APs are monitored.

Both passive and active scanning generally involve *switching* the WLAN radio to *different channels* since in real networks, APs are placed in non-overlapping frequencies (a.k.a channels) in order to minimize radio interference between them. However, in commercial and public Wi-Fi networks, APs can be found all over the 2.4GHz ISM spectrum (IEEE 802.11b/g). Thus, quantitatively, the generalized expression for handoff delay ($\mathcal{D}_{handoff}$) can be expressed as:

$$\mathcal{D}_{handoff} = T_{cs} * (N_c + 1) + N_c * \mathcal{D}_{scan} + \mathcal{D}_{auth} + \mathcal{D}_{asso} \quad (5.1)$$

where T_{cs} and N_c are the channel switching time and the number of channels, respectively. After scanning all the channels, the MN switches to the channel which has the best AP and this accounts for $(N_c + 1)$. The value of T_{cs} varies across wireless chipsets but usually lies in the range of 5ms and 19ms [65]. The scanning delay, \mathcal{D}_{scan} , has different values depending on the scan type (passive or active) while both \mathcal{D}_{auth} (authentication delay), and \mathcal{D}_{asso} (association delay) usually lie below 10 ms.

Using Equation (5.1), we can see that for *passive scan*, Layer 2 handoff delay is of the order of 1000ms for $N_c = 10$ and considering the beacon inter-arrival time to be 100ms. For the case of *active scan*, the handoff delay actually incurred depends on the maximum time the MN needs to dwell in a particular channel waiting for probe response message. Empirical studies [71] suggest that this value ranges upto 7ms. Dwelling in the channel for a higher time, however, does increase the probability of finding more APs. The usual handoff delay using active scanning has been observed to vary between 350ms to 500ms [65]. Next we introduce our proposed framework for achieving fast seamless handoff.

5.2 Dynamics of the Beacon Inter-arrival Time

In the infrastructure mode, the APs usually sends out beacon frames at fixed intervals of around 100 – 102 ms [65]. The information present in the beacon frames is used to (i) provide information about the presence of APs to MNs and, (ii) maintain the timing synchronization between the AP and the associated MNs. We, show that the inter-arrival time of the beacon arrival process in an MN is in general, a non-stationary stochastic process. Since the nature of data collected would play an important role in our analysis, let us first look into how the beacon frames are collected for our work.

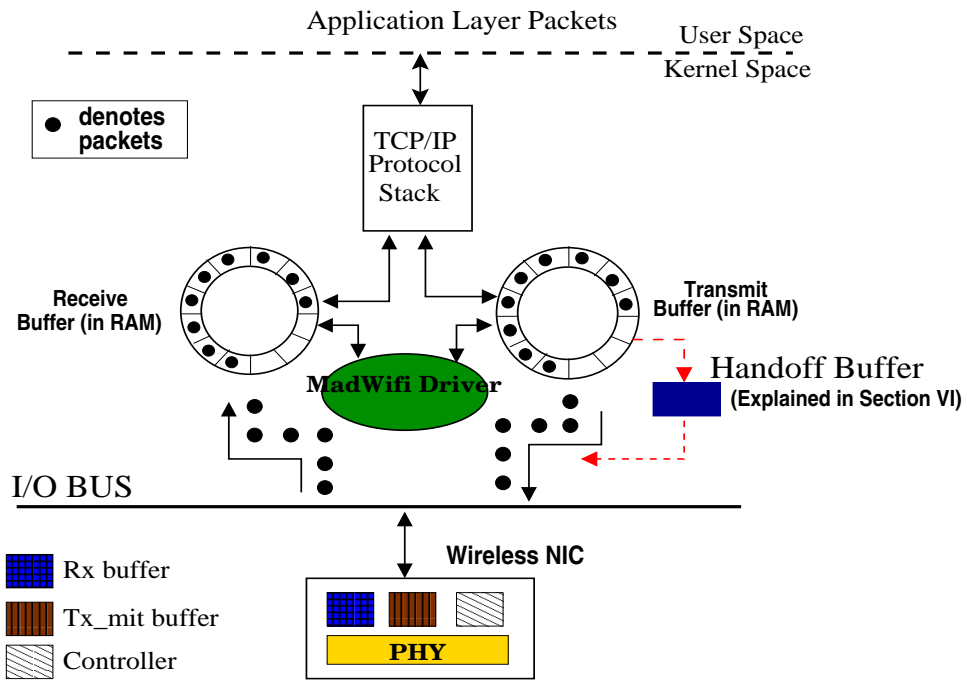


Figure 5.1. High level architecture of the proposed framework.

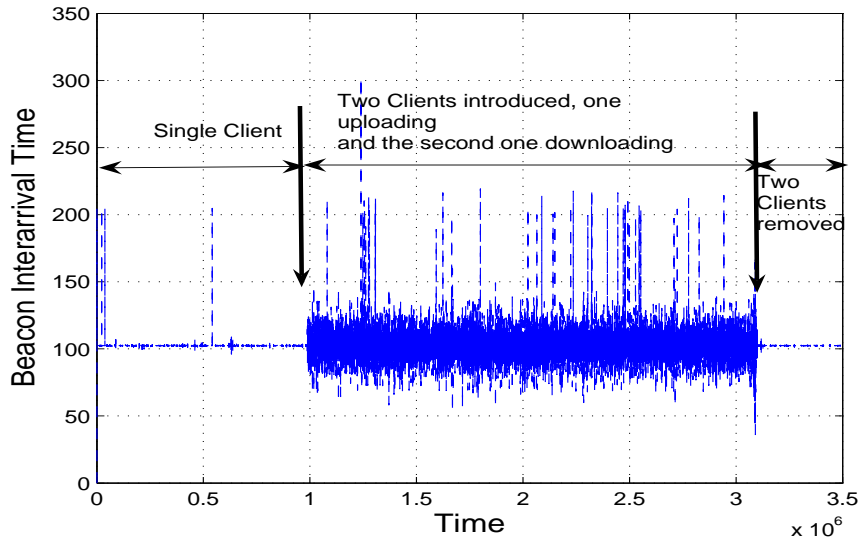


Figure 5.2. Beacon Inter-arrival Time calculated using data collected from our experimental testbed (IEEE 802.11a). As the number of clients accessing the wireless medium increases, it is observed that the beacon inter-arrival time becomes less deterministic in nature.

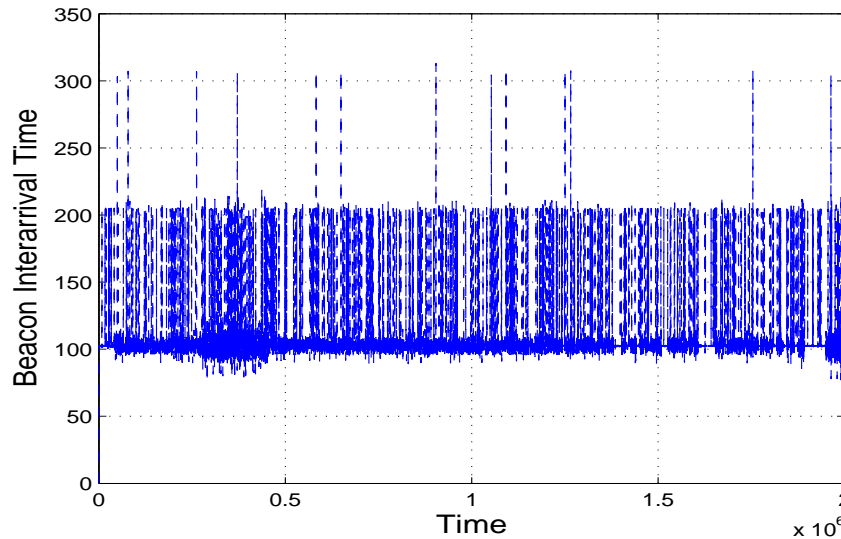


Figure 5.3. Beacon Inter-arrival Time calculated from data collected from UTA wireless network. The data was collected on a weekday from IEEE 802.11g network.

5.2.1 Empirical Data Collection and Initial Observation

We collected beacons from three different sources: (i) from our wireless testbed with six APs operating in IEEE 802.11a mode (Figure 5.2) (ii) from a IEEE 802.11g WLAN inside UTA and (iii) from a *commercial hotspot* inside Gatwick airport, U.K. This network was observed to operate in IEEE 802.11b mode. The time series of the beacon inter-arrival processes are shown in Figures 5.2, 5.3 and 5.4.

In each of the cases, the measurement was carried out by inserting probes at the `net 802.11` layer inside the MadWIFI driver [48]. We measured the inter-arrival time of the beacons from the associated AP. In the case of data collected from our wireless testbed, experiments were performed at 5.18 GHz so as to reduce interference from existing deployed 802.11b/g networks. The measurement was first carried out with one AP and a single MN. After about 15 minutes, two MNs were introduced in the same network and forced to get associated with the same AP as the first MN and stay connected for about 30 minutes. All the MNs independently used *sftp* for uploading/downloading a

10MB file from the multimedia server (refer Figure 5.2).

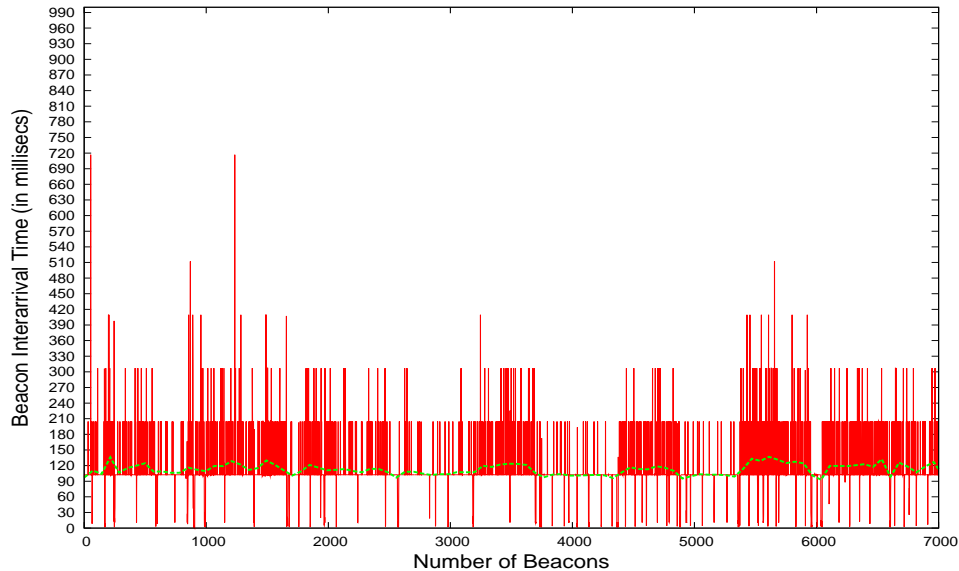


Figure 5.4. Beacon Inter-arrival Time calculated from data collected from a commercial hotspot at Gatwick Airport, UK, April 2005. Network was in IEEE 802.11b mode.

This simple experiment provides us with an insight that *probably* in the absence (or low volume) of background traffic and competing channel access, the beacon inter-arrival is deterministic in nature with values closely equal to the beacon generation inter-arrival time. However, the beacon inter-arrival time starts to become random as soon as more than one MNs starts sharing the wireless channel.

Further measurements carried out with live networks involving the UTA academic campus network and a commercial AP at Gatwick Airport, UK corroborate our initial findings. But what is the *underlying network behavior* that is causing the beacon process to become non-deterministic? We provide answers to such a question in the next section of our work.

5.2.2 Why is the Beacon Process Non-Deterministic?

In order to investigate the underlying process of beacon arrival and generation, we need to clearly understand the path of beacon traversal as it is released from the AP and arrives at the MN. In Figure 5.1, we presented the architectural view of a WNIC. The wireless driver, (MadWifi in our case), embedded inside the OS kernel (Linux in our case) contains the *device independent* 802.11 state machine.

Process of packet traversal: Management and control frames are generated by the Madwifi driver based on the IEEE 802.11 state machine present in the `net802.11` layer while the data packets are generated by the TCP/IP stack (from application layer). The details of the state machine is explained in Section 5.3.1. While *transmitting*, the driver attaches *higher priority* to management frames over the data packets and maintains *two different queues* inside the WNIC as shown in Figure 5.1. However, only *one queue* is maintained while *receiving* packets/frames from the WNIC. The receiving frame types distinguished only at the `net802.11` layer present inside the driver.

Process of beacon generation at the AP: Following the IEEE 802.11 standard [3] requirements, a timer inside the wireless chipset generates a *hardware interrupt* indicating the driver to create a beacon frame. The hardware interrupt fires at fixed intervals of time (typically the interval is set between 100 – 102ms). What happens aftermath is shown in Figure 5.5.

At the time instant the beacon interrupt has been raised, there *already exists* interrupts from network packets[†] which are ahead in the interrupt queue. Thus, in spite of having higher priority than the data packets, the *beacon interrupt has to wait* (queuing delay) till all the packets (or interrupts raised by the packets) have been processed by the interrupt handler. However, once the beacon interrupt has been delivered and processed

[†]In Wi-Fi networks, packets are pulled from the hardware queue using per-packet interrupts. Elegant schemes, like *interrupt coalescing*, commonly found in high speed wired networks, would be an overkill since packet inter-arrival rate is very low in current Wi-Fi networks.

by the interrupt handler, the beacon frame is immediately transmitted by the driver to the WNIC hardware from where it is sent out in the wireless medium.

If we consider the individual clients in the network as *ON-OFF* traffic sources, then the input to the onchip hardware buffer has dynamics similar to the one created when a large number of discrete *ON-OFF* are multiplexed together. The dynamics of the onchip buffer can be effectively captured by the Benes fluid queue analysis [15]. Also, it is well known that such multiplexing processes give rise to long range dependency in network traffic [39]. All these causes the *beacon generation process to be stochastic in nature*. Now let us look at the process of beacon reception at the MN.

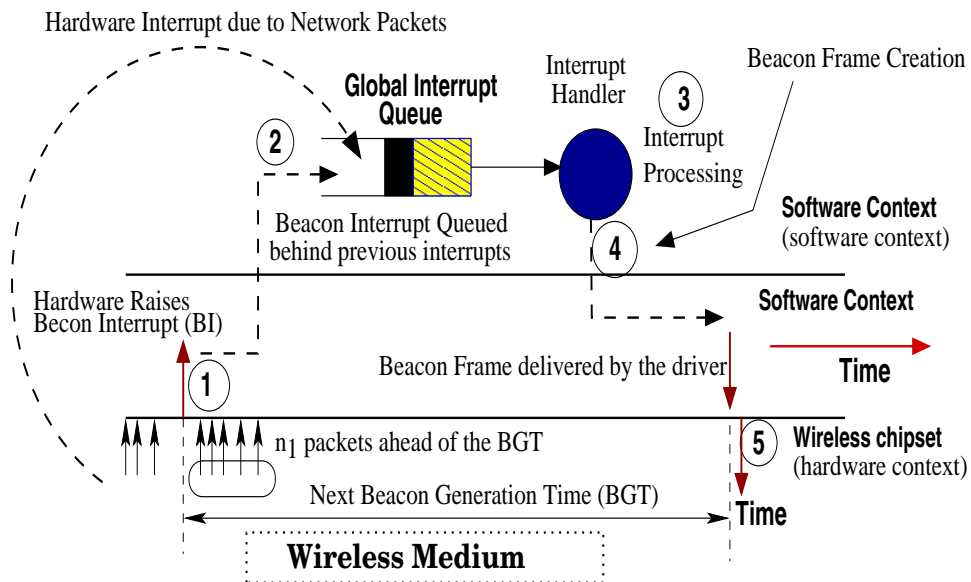


Figure 5.5. Dynamics of the beacon processing at the AP. It consist of five main steps: (1) hardware timer raises the beacon interrupt (2) the interrupt is queued in the global interrupt queue (3) interrupt handler processes the interrupt (4) the beacon frame is created inside the driver (5) the frame beacon is transmitted by the wireless chipset.

Process of beacon reception at the MN: Let us refer to Figure 5.6 and understand the beacon arrival process at the MN. Observe that the 802.11 state machine is aware that

a frame is of type beacon only after it has reached the wireless driver. Packets arriving from the wireless medium are initially placed on the onchip receive buffer of the WNIC[†].

Notice that in this case, the beacon has been queued even before it has been able to raise an interrupt. Thus, unlike the beacon generation process where the main delay is at the global interrupt queue, the beacon arrival process at the MN experiences two sources of delay - at the onchip receive buffer and at the global interrupt queue. Only after the packets ahead of the beacon in the buffer has been processed, the beacon frame is delivered to the *net802.11 layer*. Proceeding along similar lines of reasoning as in the beacon generation case, the beacon deliver process to the wireless driver can be shown to exhibit stochastic nature. However, be note that in the Atheros AR5212 chipset, the hardware exports the timestamp at the instant the beacon has arrived at the chipset. As of date, such hardware timestamp information is not used for beacon inter-arrival processing. Thus to conclude, the beacon inter-arrival process suffers from *random delay* both during transmission from the AP and as well as when being received by the MN. This inherently makes the beacon inter-arrival time non-deterministic. *Consequently, the random nature of the beacon inter-arrival time jeopardizes any handoff scheme based on deterministic nature of the beacon interval.*

5.2.3 Beacon Inter-arrival Time Sequence: Short Range Dependent

Having observed that the beacon arrival process is stochastic, let us investigate how it impacts the beacon inter-arrival time dynamics. That would eventually drive our handoff algorithm. At a first step, such behaviour can be investigated by studying the autocorrelation function (ACF) of the beacon inter-arrival time sequence.

[†]Only packets whose MAC addresses match the WNIC address are passed to the wireless driver, assuming that the driver is not working in promiscuous mode.

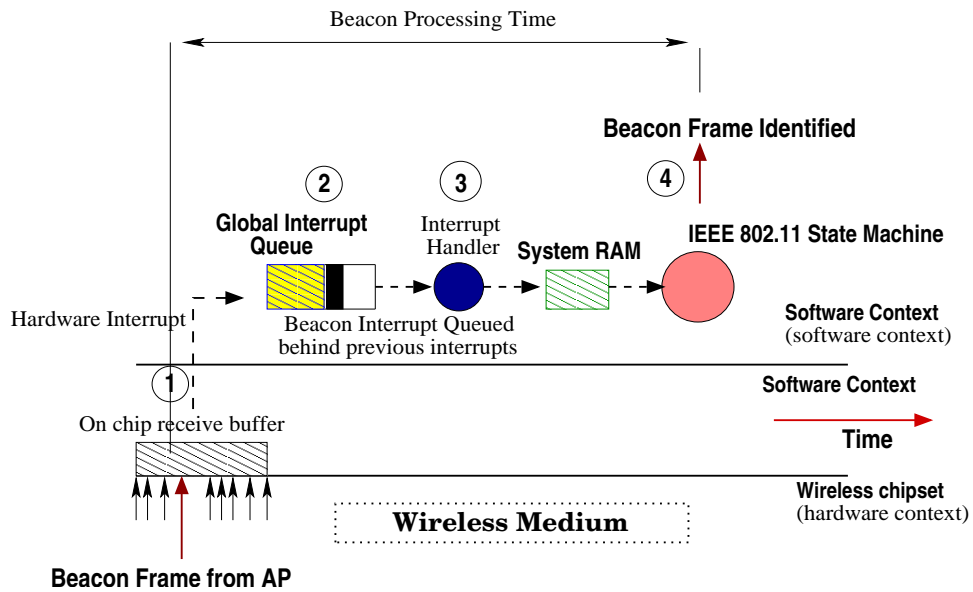


Figure 5.6. Dynamics of the beacon processing at the Mobile Node (MN). It consists of four main steps: (1) beacons arrive from the wireless medium and are kept on onchip receive buffer. The hardware raises an interrupt (2) the interrupt is queued in the global interrupt queue; (3) interrupt handler processes the interrupt and is moved to the system RAM. (4) the beacon frame enters the IEEE 802.11 state machine and is identified.

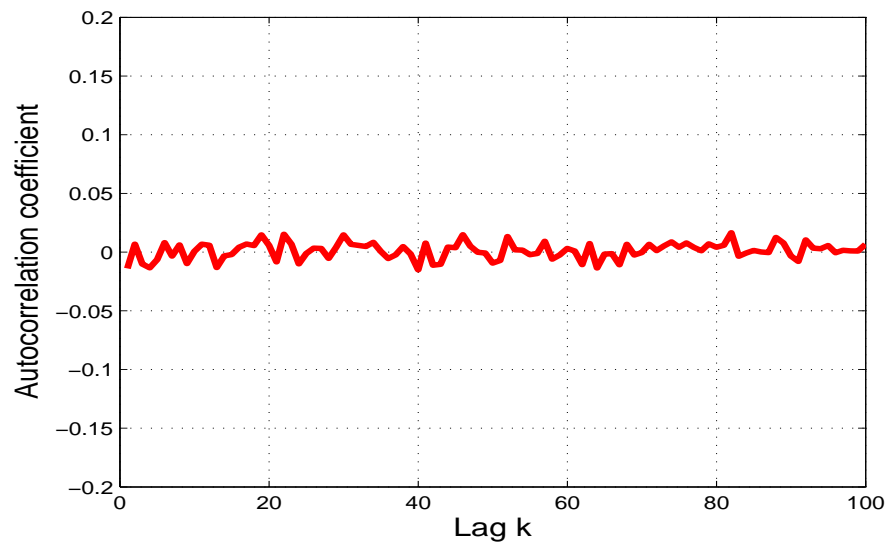


Figure 5.7. Autocorrelation structure in beacon interarrival.

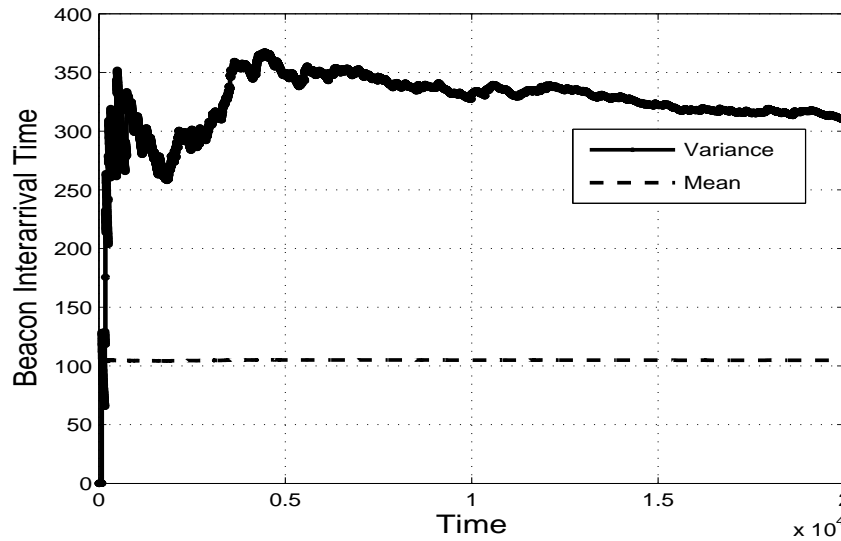


Figure 5.8. Variation of mean and variance of the beacon interarrival sequence at the MN.

The autocorrelation function (ACF), $r(k)$, which is a measure of similarity between the sequence $x_n(t)$ and a time shifted version of itself, $x_n(t+k)$, is given by:

$$r(k) = \frac{E[(x_t - \mu)(x_{t+k} - \mu)]}{\sigma^2}$$

Figures 5.7 and 5.8, plot the ACF w.r.t lag coefficient k and the mean variance time sequence. We observe *very little correlations* between distant samples indicating that the beacon inter-arrival sequence is short range dependent (SRD). As a matter of fact, the series is wide sense stationary (constant mean, finite variance). Without going into the details, it can be shown that the beacon inter-arrival time can be mapped to a *shot noise processes*.

5.3 Framework for fast Handoff

The proposed framework constitutes of a suite of algorithms that work in close tandem for handling the Layer2 handoff process along with the associated features of dynamic discovery and monitoring of APs and maintenance of the handoff buffer. We have already seen how the AP scanning phase is the bottleneck for fast handoff and how an MN might be able to reduce handoff delays provided the list of accessible APs is known apriori. However, maintaining such a list of APs has the following challenges: (i) the number of APs in the backup and the criterion for selecting the APs; (ii) frequency of monitoring the APs in the backup list; (iii) criterion for discarding the APs from backup list; (iv) frequency of triggering the service discovery mechanism ; and (v) simultaneously monitoring the current associated AP so as to determine when to actually trigger the handoff. In the following, we explain how the above challenges are addressed and how actually a handoff is executed in 802.11 protocol.

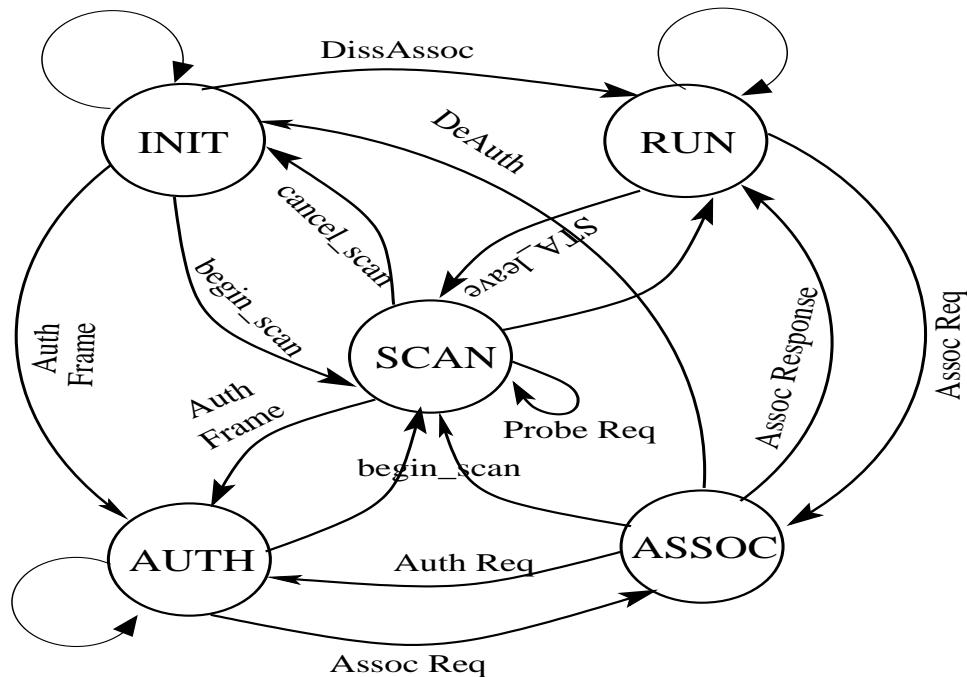


Figure 5.9. IEEE 802.11 Finite State Machine (FSM) present inside the wireless driver.

5.3.1 IEEE 802.11 State Machine

The IEEE 802.11 FSM shown in Figure 5.9 is what primarily drives the WNIC. This entire state machine implemented as `net802.11` layer inside the driver is independent of any underlying wireless hardware and is common across almost all IEEE 802.11 devices. Major portions of our framework has been implemented in this `net80211` layer with minor changes to the available chipset parameters (Atheros AR5212 in our case). Hence, with minor modifications for the chipset portion, our algorithms should be portable across all IEEE 802.11 compliant devices. The heart of the 802.11 network is the state machine shown in Figure 5.9 which drives the operation of the WNIC card. This entire state machine has been implemented as a *net80211 layer* in the MadWiFi driver. The *net80211 layer* is device-independent and hence should be able to work with any other 802.11 device. Our proposed solution has been mostly implemented in the *net802.11 layer* with minor changes to the atheros specific chipset constraints. Hence with minor modifications, our approach should work for any 802.11 device. There exist 5 states for the 802.11 state machine : (i) **INIT**, (ii) **SCAN**, (iii) **AUTH**, (iv) **ASSOC** and (v) **RUN**. The valid transitions between the states are shown in Figure 5.9. The following state transitions are considered invalid: $\langle AUTH \rightarrow INIT, RUN \rangle$, $\langle INIT, SCAN, ASSOC \rightarrow ASSOC \rangle$. On bootup, the card is driven into the **INIT** state from where the 802.11 protocol drives it to the **SCAN** mode. Active scanning of the entire spectrum for determining the best possible AP to which it can get associated with is enforced initially in the **SCAN** state.

RSSI as an Indicator of Link Conditions: The RF energy measured by the radio inside the wireless chipset is reported as a numeric value. This value, referred to as the Received Signal Strength Indicator (RSSI) is frequently used as a metric for measuring the wireless link conditions. While higher values indicate better link conditions, the range of values reported by the WNIC varies across chipset vendors. The RSSI values reported by AR5212 chipset varies from 0 to 60. The RSSI measurements obtained cannot be

directly mapped with the Signal to Noise Ratio (SNR) of the wireless channel due to the nonlinear behavior of the RSSI detector. However, every chipset provider provide a conversion table between RSSI and SNR. In our approach, we have used weighted SNR (computed from the RSSI values) of past history (Exponential Weighted Moving Average technique) to get an estimate of the link conditions of the APs.

After scanning, the FSM moves to **AUTH** state where it needs to authenticate itself with the AP with which it desires to be associated with. The management frame **Assoc Req** drives the FSM to **ASSOC** state. Data transmission is initiated by the driver only when the FSM is in the **RUN** state after receiving the **ASSOC Resp** frame from the AP. As of now similar to existing solutions, our solution handles seamless mobility for open authentication APs only. The best possible AP is the one whose Received Signal Strength Indication (RSSI) has the highest value. Note that the RSSI value is exported to the MadWiFi driver by the chipset. Unfortunately, one random RSSI value is only what the 802.11 device currently utilizes to determine the initial best AP. On full scan completion and AP determination, the state machine moves to **AUTH** state where basically it authenticates itself to the AP to which it desires to be Associated with. So far we have implemented out seamless mobility solution for open authentication APs only[†]. The MN on receiving the **Assoc Req** frame moves to the **Assoc** state where it actually completes all the management issues of getting associated with the AP. The data transmission is initiated only when the 802.11 state machine receives the **ASSOC Resp** frame and drives the state to the **RUN** state. Having explained the 802.11 protocol states briefly, we proceed to present in details how our approach works without modifying the native IEEE 802.11 protocol introduces scanning, both passive and active leading to automatic AP discovery, as well as seamless handoff.

[†]We plan to implement Authenticated seamless connection for our next release

5.3.2 Bootup Phase and Creation of Initial Backup AP List

Our algorithm is triggered at boot time when the wireless driver is loaded in the kernel and the WNIC detected. The Atheros AR5212 chipset starts scanning with channel 1 (2.412) GHz and updates the node table (*ni_table*) when it receives probe response packet. Under the circumstances if the client was not able to detect a single AP, the entire process is repeated. All the APs which respond at this stage are likely candidates with which the MN can get associated with. During this phase, our algorithm does not interfere with the existing normal operations of the WNIC but monitors the list of all the APs which the WNIC has been able to detect. We term this phase as the *information gathering phase*. The relevant information stored for each AP are the frequency/channel where it resides, the MAC address, Extended Service Set Identification (ESSID) and corresponding RSSI values. It should be noted that this list is *not* our backup AP list but mainly serves to indicate the number of APs the MN can listen to. On completion of the full scan and after the WNIC has been associated in open authentication mode by the driver, we sort the channel list based on the frequency of APs present and determines a set of k (set to 2 in our experiments) APs based on their respective RSSI values. This is the initial backup list. We closely monitor the APs in this list for signal strength variations. The k APs need not necessarily be on the same channel. It should also be noted that the AP with the highest RSSI value is not present in the backup list since it is usually the AP that the WNIC has already been associated with.

5.3.3 Service Discovery and Maintenance of APs

Background channel scanning is extensively used mainly for two purposes: (i) to *monitor* the backup AP list, and (ii) to *create* the backup AP list during the initial stage or when the list is observed to be empty or depleted. It should be noted here that maintaining a large backup list does enhance the reliability in terms of having a long

list standby APs at the cost of increased overhead of AP monitoring. In the worst case scenario, when all the backup APs are in *different channels*, then the penalty of having a large backup list becomes evident. Such overhead comes from the fact that the WNIC has to switch the radio to the desired channel and also switch back to the channel in which it currently resides. During this transition phase, no data packet can be sent out and the throughput is penalized.

AP Maintenance: In the maintenance phase, we sequentially monitor each of the APs in the backup list. This monitoring mechanism is driven by the asynchronous events of (i) the number of APs present in the backup list at any instant of time, and (ii) the weighted average RSSI value of each of the APs. We have considered an Exponential Weighted Moving Average (EWMA) of the current and the last w RSSI samples for each AP being monitored. In our experiments, we noticed that setting the value of w to more than 20 does not yield any additional benefits. The EWMA technique filters out jitters and provides a true status of the the signal strength of the AP. In our implementation, we have used Linux *kernel timers* to implement this process as presented in the following code snippet:

```
mod_timer(&ic → ieee80211_passive_scan_timer,
          jiffies + (timeToFire*HZ)/1000);
```

We initially set the `timeToFire` variable to a timeout of 20 ms after the timer has been fired (for AR5212 chipset). This is due to the fact that probe responses usually have an average delay of 7ms and the channel switching time is of the order 10 ms in AR 5212 chipset. Prior to switching the channel, it is usual to stop the the hardware and driver queues. However, we use packet buffering to provide the illusion of the interface being always “ON”. It also sends power saving frames so that packets get buffered at the AP

during channel switching.

On switching to the desired channel, a probe request message is broadcasted. If the probe response message is received successfully, the average RSSI value of the AP is updated and it is marked “alive”. After all the backup APs has been successfully monitored, the backup list is reordered on the basis of the average RSSI value. This ensures that the handoff always takes place with the best available AP.

Service Discovery: The Service discovery mechanism has been designed to be *adaptive* w.r.t the existing number of backup APs and is usually triggered (i) during the bootup phase and (ii) when the backup AP list becomes empty. Once the backup list of APs gets populated, then the maintenance and discovery mechanisms run in tandem. The discovery mechanism, however, stops once the backup AP list has been populated with the desired number of APs.

The service discovery mechanism first selects the list of channels on which the backup APs are present and then chooses the channel from this list on the basis of non-overlapping criterion between both the existing AP and backup APs. Next we describe how the handoff is actually executed.

5.3.4 Handoff Execution

In order to execute a successful handoff decision, the algorithm needs to know to which AP it needs to switch to and how to identify the time instant when the handoff process should be triggered. We continuously monitor the RSSI value of the associated AP using EWMA technique. The decreasing RSSI trend signifies the fading of the signal strength and possibility of a handoff. Handoff is triggered when the effective RSSI value goes below a specific threshold that depends on the wireless radio being used. In AR5212 chipset, an RSSI value of less than 18 is considered to indicate low link conditions. The usage of EWMA technique for effective RSSI computation eliminates the ping pong ef-

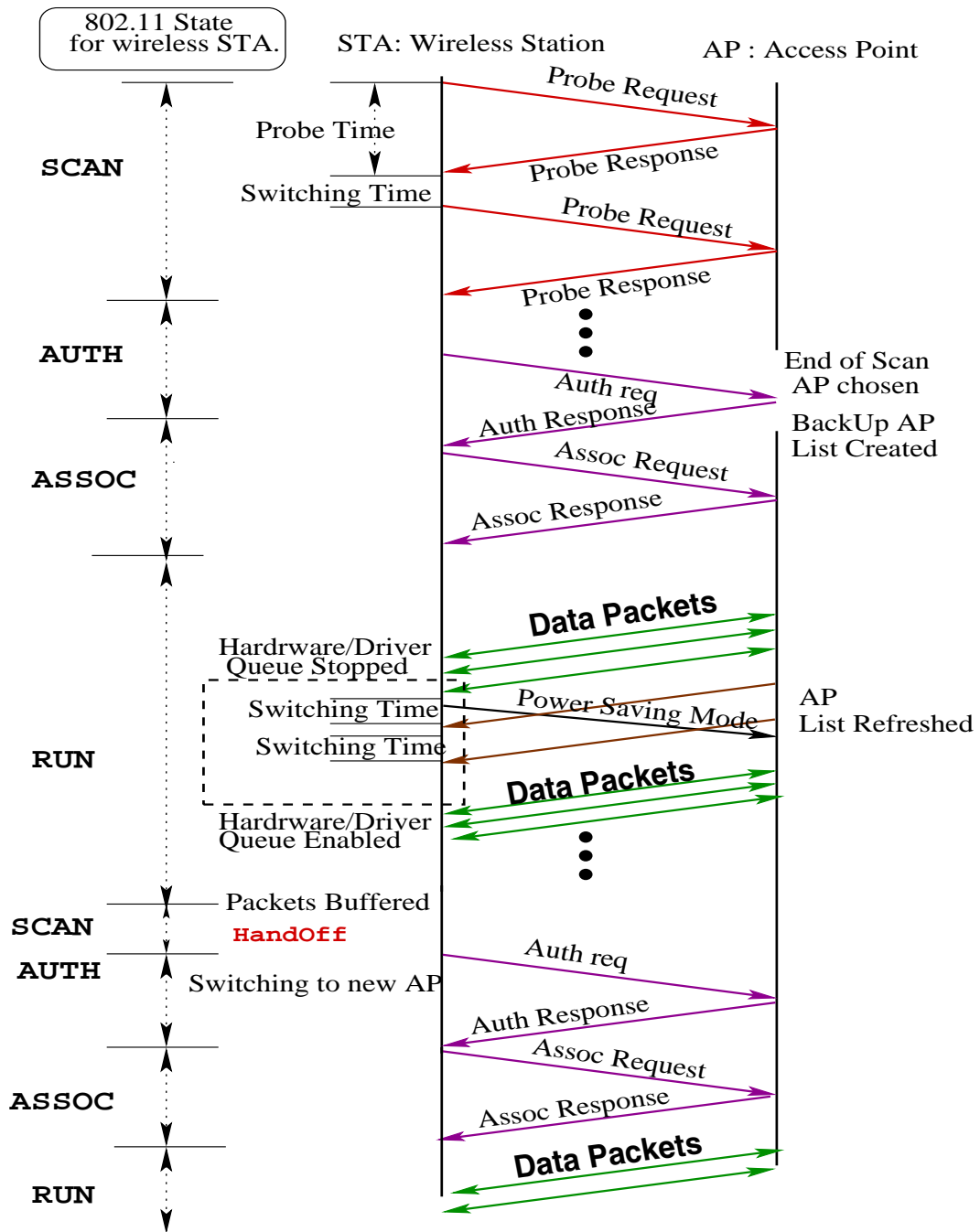


Figure 5.10. Communication diagram showing the message flow between our proposed framework and the Wi-Fi network.

fect.

What are the steps involved? The actual handoff execution encompasses the follow-

ing six steps: (i) the hardware and driver queues are stopped to prevent data loss and a disassociation message is sent to the currently associated AP with the help of packet buffering that depends on the rate of data traffic; (ii) the state of the FSM is switched from the **RUN** state to the **SCAN** state. (iii) the next AP to which the MN will get associated with is selected from the backup AP list. That is is already ordered according to the average RSSI values; (iv) the WNIC is switched to the specific channel and it fires an **Assoc Req** message; (v) the state machine is switched to the **AUTH** state where it dwells till it receives an authentication response; and (vi) on successful association, data packets are sent out by the driver.

The messaging diagram in Figure 5.10 illustrates the entire handoff process. Note that direct transition from **RUN** to **AUTH** state breaks the IEEE 802.11 state machine. Hence, we transition through the **SCAN** state though in practice we do not enter the traditional “scanning” phase. Consequently, no time is wasted in the FSM **SCAN** state while executing the handoff. During this entire procedure the solution ensures that all the maintainance, monitoring and service discovery functionality for the backup AP list is paused.

5.4 Performance Evaluation

The proposed solution for seamless fast handoff in Wi-Fi networks has been implemented inside the MadWiFi driver as part of the *net802.11/wlan layer* and with minor modifications to the *ath layer(if_ath.c)* [48]. The MadWiFi driver works for the widely used Atheros [7] based WNIC chipsets (AR5212,AR5211, AR5210). The version details of the MadWiFi driver on which our algorithm has been implemented is as follows : wlan (v0.8.6.0), ath_pci (v0.9.6.0) and ath_hal(v0.9.14.9).

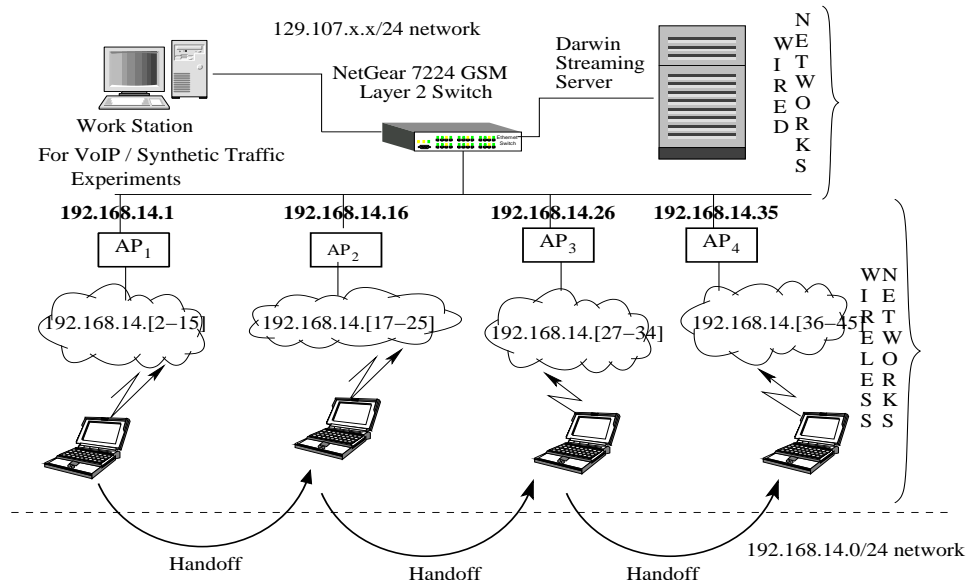


Figure 5.11. Illustration of the handoff process. This is also the experimental testbed we have used for evaluating our proposed solution. AP_1 , AP_2 , AP_3 and AP_4 shown are a combination of commercial APs (Linksys WRT54G) and standard Linux PCs running MadWifi driver.

5.4.1 Experimental Testbed

A Sharp Actius (PC-MP30) laptop with 512MB RAM and an externally supplied 32bit cardbus from Netgear (WAG511v2) is used as the mobile client. We make sure that during the entire experiment the onboard wireless chipset is turned off. The development platforms (including the mobile laptop) run the same configuration of SuSE Linux with unpatched vanilla 2.6.11.6 kernel. The experimental testbed used to evaluate the performance of our algorithm is shown in Figure 5.11. It comprises of two Linksys (model WRT54G) commercial APs and two standard Linux PCs with 32bit PCI WNIC (Netgear WAG311) running the MadWifi driver in the master mode. Throughout the experiments we have considered two cases: in the first case the backup APs and the current AP are overlapping channels. The second case the backup APs are in non-overlapping APs. The testbed comprises of two subnets: a Gigabit Ethernet wired subnet (129.107.x.x) and an IEEE 802.11g wireless subnet (192.168.14.x). The content delivery server running

the Darwin Streaming Server [23] (a.b.c.114) is connected to the workstation (a.b.c.228) through a Layer 2 managed Gigabit Ethernet switch (Netgear GSM 7224). The APs were separated by more than 60 feet. This distance was sufficient enough such that the average RSSI went below the threshold and handoff was automatically triggered.

5.4.2 Experimental Results

Base Case Comparison: Ideally speaking, true performance gain achieved can be measured if we compare it with the scheme when no handoff algorithm exists. However, both in windows and Linux systems, complete connection is lost and full scan is enforced and at times manual intervention is needed rendering accurate measurement infeasible. In the base case, the observed handoff delay is between 1 to 2 minutes. On the average, for our proposed scheme the handoff delay varies between 15 to 20 ms respectively and 20 to 26 ms for overlapping and non-overlapping channels.

Performance Evaluation With VoIP and Multimedia Streaming: We evaluate the performance of our proposed handoff scheme with respect to different application requirements.

First, we employ *ping* to generate ICMP echo request packets at periodic intervals varying between 20ms to 100ms to mimic VoIP and multimedia stream requirements respectively. Thereafter, we generated actual VoIP and multimedia streams with *Skype* and *Mplayer* [55] respectively. The packet level statistics collected during each run of the experiment provide valuable insights about the impact of our algorithm on application level performance. The metrics chosen for evaluation are: round trip time (RTT), packet inter-arrival time (IAT), and percentage of packet loss. As a measurement tool, we used *Ethereal*(v0.10.12) for capturing the packet level statistics at the MN. For VoIP experiments (based on *Skype*), we wrote our own analysis tool using the *libpcap* packet library for extracting the necessary statistics. In all the experiments that we have conducted,

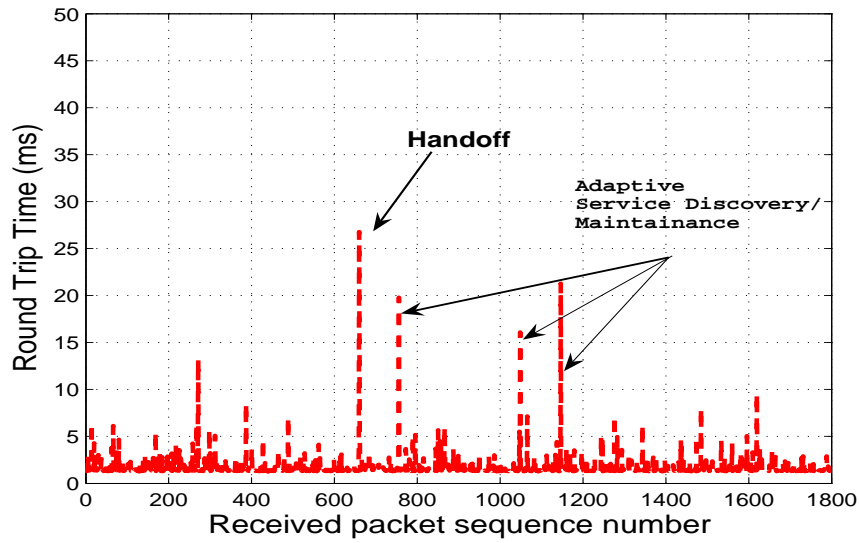


Figure 5.12. Graph of Round Trip Time (RTT) vs received packet for packets generated at 20ms .

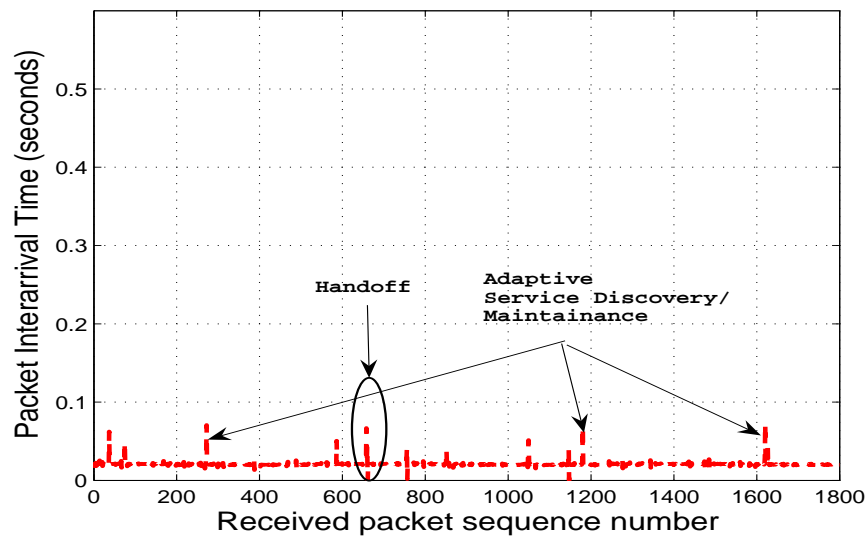


Figure 5.13. Graph of Inter-Arrival Time (IAT) vs received packet for packets generated at 20ms .

the delay suffered during handoff for different traffic types using our handoff solution varies between 15ms to 26ms.

In Figures 5.12, 5.13, we illustrate the performance of our solution when ICMP

Echo Requests packets at 20ms interval was generated from the MN to a server in the wired domain. The backup APs were operating in non-overlapping channels. We plot the RTT versus the received packet train as the MN moved between different APs. It is observed from Figure 5.12 that the handoff delay using our algorithm is of the order of 26ms with no packet loss. The other spikes in Figure 5.12 signify the backup list monitoring and maintenance operation. This is reflected by the increase in the RTT as shown in Figure 5.12 at specific instances when the WNIC switched channels and performed service discovery or monitoring. We notice that at certain instances the packet inter-arrival drops below the average interarrival time. This is because we *buffer the packets* by stopping the queue (or enabling the handoff buffer) before triggering a handoff or a maintenance/service discovery phase. This results in back-to-back packets which reduce the inter-arrival time. It is to be noted that *no packet loss* was observed during the handoff.

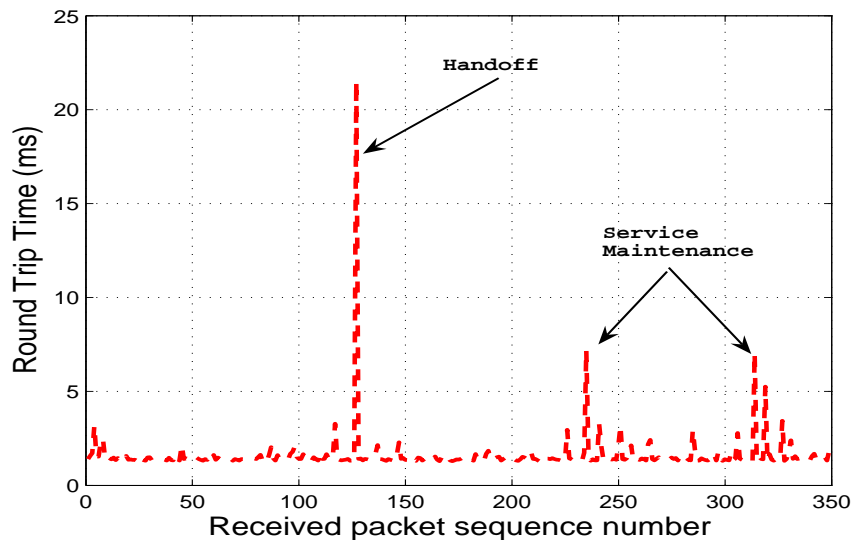


Figure 5.14. Graph of RTT versus received packet for packets generated at 100ms .

Similar experiments were conducted by setting the ping packet generation rate to 100ms and the backup APs were all in overlapping channels. The results are shown in Figures 5.14 and 5.15. The same trends as above are observed in these experiments also. However, the service discovery mechanism is triggered at *different instances* for different experiments since it depends on the *dynamic status* of the backup AP list. The handoff delay is reduced to 15ms on account of the backup AP being in the overlapping channel. Similarly, the monitoring time is also reduced in comparison to the previous case where the backup APs resided in non-overlapping channels. VoIP calls were made from the

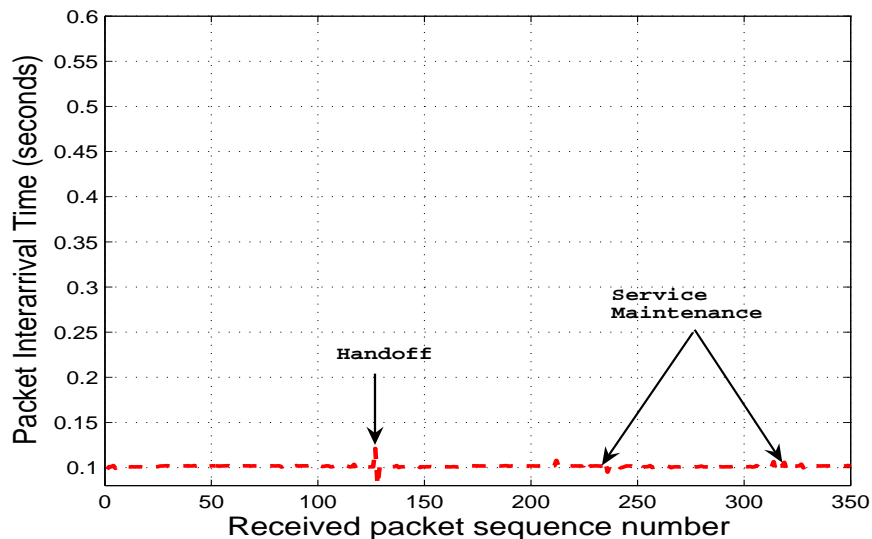


Figure 5.15. Graph of packet IAT vs. received packet for packets generated at 100ms.

MN using Skype to a server in the wired domain. On an average, we have observed that delay for the inter-arrival time is bounded by 26 ms and we present experimental data in Figure 5.16 for VoIP experiments. Experiments involving Darwin Streaming Server streaming 64Kbps MPEG-4 bitstream to the MN using RTSP protocol was also conducted. *We have observed that the handoff delay using our solution does not affect the*

performance of the inter-arrival time of the frames and is again bounded by 26ms. As seen in Figure 5.17, the handoff takes place between packets 1500 and 2000, but there is no perceptible change in the mean packet inter-arrival time. To qualitatively evaluate, we conducted two experiments - one in which the backup APs were in non-overlapping channels and in the other case they were in overlapping channels. We streamed video and enforced handoffs in both the cases and dumped the stream to catch the frames. In case of overlapping channels, it is interesting to note in Figures 5.18 and 5.19 that no visible differences exist between the original frame and the frame generated during the handoff.

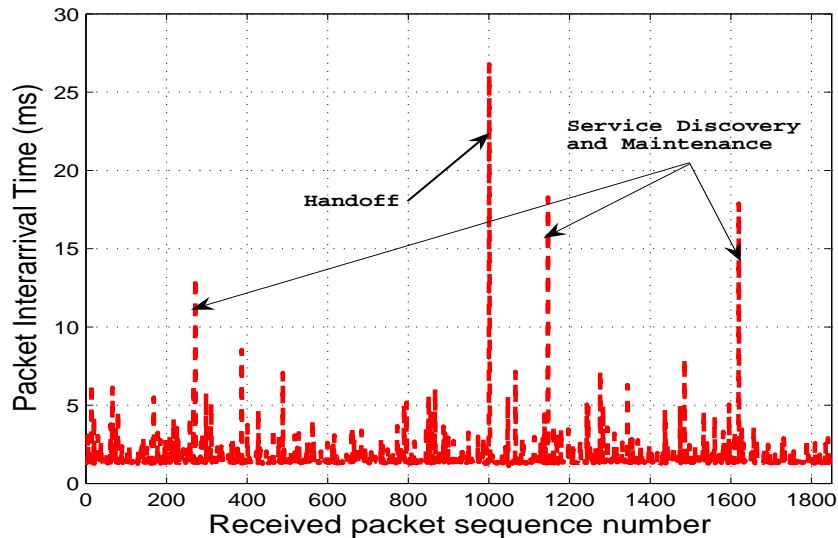


Figure 5.16. Packet IAT for VoIP stream as experienced during handoff using proposed solution.

5.4.3 Impact of Buffering

In order to ascertain the utility of the inline packet buffering mechanism, we conducted experiments with and without the buffer enabled. We used the ping application

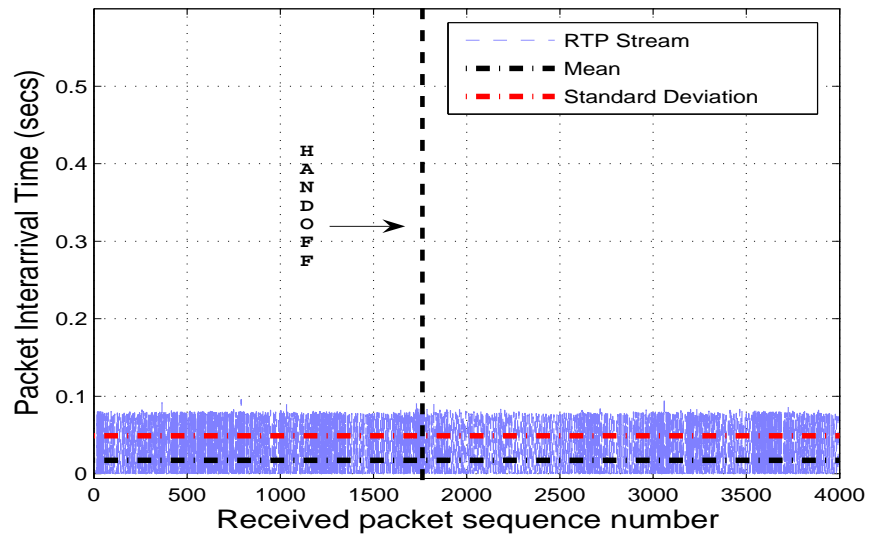


Figure 5.17. Packet IAT time for Video Stream as experienced during handoff proposed solution.



Figure 5.18. Original Image quality.

to generate packets with inter-arrival time of 5 ms, 50 ms and 200ms from the MN to the work station. Without the buffering we observe that packet loss is high as 70% whereas



Figure 5.19. Image quality with Handoff.

with the buffer enabled it is significantly lower. Figure 5.20 presents how buffering reduces the packet loss significantly.

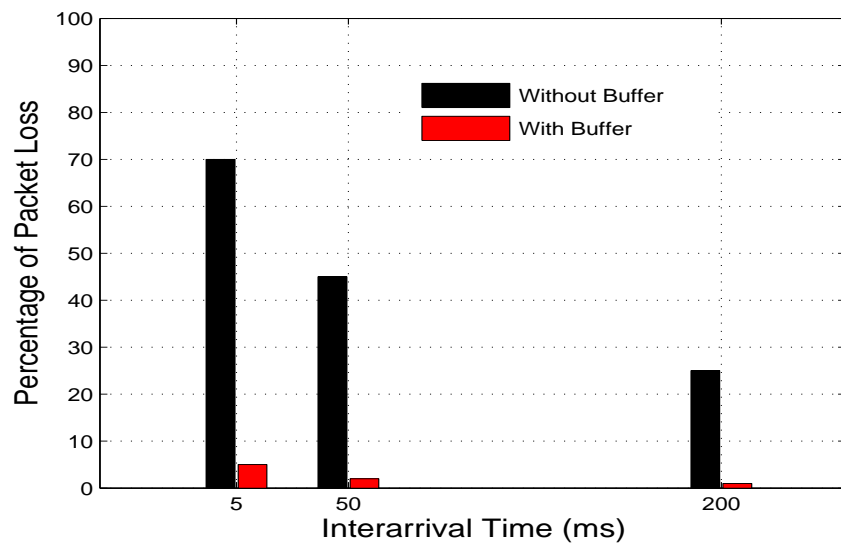


Figure 5.20. Effect of Buffering on Packet Loss for different arrival rates.

5.5 Comparison with Existing Approaches

A number of ways to reduce handoff latency in 802.11 networks have been proposed in recent years [70, 71, 65, 19, 78]. In [62], the authors proposed a *reduced channel scanning method* based on the location and distribution of APs in the network. In such a scheme, the *network has to explicitly provide support* for messages that would enable an MN to discover APs in close proximity. Similar topology based *neighbor graphs* has been put forward in [70]. In [58], the authors focused on reducing the authentication latency during handoff based on user’s registration patterns. As also observed in [65], the above solutions are (i) not within the framework of IEEE 802.11 protocol and (ii) require extensive network support. Hence, they have limited impact in providing practical handoff solutions. In [19], the authors have suggested using a *dual radio* WNIC so

Table 5.1. Comparison of Proposed Solution with existing WLAN handoff solutions

Feature	Neighbor Graph [70]	Synscan [65]	MultiScan [19]	Smart Trigger Scheme [53]	Proposed Scheme
Handoff delay (ms)	> 40	Specifically Not Mentioned (Imperceptible)	0	30 – 40	15 – 26
Packet Loss	Yes	Yes	No	Yes	Yes (< 1%)
Changes to IEEE 802.11 protocol	Yes	No (APs need NTP Synchron.)	No (Need 2 WNICs)	No	No
Software Modification ?	Yes (both at AP and Client End)	Yes (both at AP and Client End)	Yes (at client end)	Yes (at client end)	Yes (at client end)
Data Buffering during handoff	No	Yes	No	No	Yes
Automatic Service Discovery	No	No	No	Yes (only overlapping channels)	Yes
Scheme Tested with live network	No	Yes (software available)	No	No (as of yet not available)	Yes (software available)

that one radio can be used for background scanning while the other radio can be used for packet transmission. The authors even claim of achieving zero seconds handoff latency which is extremely difficult to achieve. Unfortunately the paper has used only “kernel level simulation” results.

One of the earlier and practical solutions requiring minimal network upgrade is SyncScan [65]. It is accompanied with the caveat that for fast handoff decisions, Syn-

Scan requires all the MNs and APs to be time synchronized via some time synchronizing protocol like Network Time Protocol (NTP). Recently, smart algorithms to trigger hand-off in 802.11 systems have been developed in [53]. The authors successfully used the leakage power to capture information of APs located in the overlapping frequency. Next, we compare our proposed approach with the existing ones.

Comparison of 802.11 Handoff Solutions : In Table II, we compare both quantitatively and qualitatively our scheme with the different handoff solutions based on the desirable features for a 802.11 handoff solution. It is to be noted that the current state of the art 802.11 networks do not have any sophisticated handoff scheme and suffer handoff delays of the order of seconds; hence we did not put it in the comparison table. Our proposed solution achieves the least handoff delay in comparison with the other existing solutions. Moreover, on account of its inline packet buffering feature it achieves the lowest packet loss. The most important criterion in which our scheme and the one in [53] excel is that neither of them requires any change to the IEEE 802.11 protocol or network infrastructure support. Furthermore, our scheme surpasses [53] on account of its superior AP discovery mechanism encompassing the entire spectrum.

Limitations of Our Scheme and Future Work : The major limitation of proposed scheme is that the frequency of the monitoring the backup APs and the number of APs in the backup list is adhoc. Also during background monitoring, though our scheme does not loose data packets but introduces some delay. However, as observed from the experiments we have seen the delay is less than 15ms and is acceptable for real time applications. We plan to analytically formulate and tune the frequency of background monitoring and also the optimum number of backup APs. In addition, we would like to make our scheme successful for authenticated APs and work with IAPP in tandem.

5.6 Summary

With the rapid deployment of Wi-Fi networks, fast roaming will increasingly become important. Aggressive real time applications like VoIP demand latency less than 50 ms. In this chapter, we have presented the first *practical* and *available client side* solution framework that is capable of reducing handoff delays in Wi-Fi networks within acceptable limits. It works with a IEEE 802.11 compliant WNIC card with a *single radio*, does *not* require any changes to the network infrastructure and best of all is available as a driver update to the client. Our proposed solution introduces the following features: (i) *dynamic* and *adaptive* AP discovery mechanism in *realtime*, (ii) *inline packet buffering* wherein the operating system (OS) stack is made oblivious to the handoff process and (iii) *software leakage filter* using which stray beacons from different channels are effectively removed. We have successfully implemented and tested our algorithm on Atheros AR5212 chipsets using the MadWifi driver and evaluated its performance using different traffic classes in both commercial networks and controlled wireless testbed. In all the cases the handoff performance was bounded by 26ms at the worst case and 20ms on the average.

CHAPTER 6

CONCLUSIONS

The dissertation broadly addresses the issues pertaining to data services in multi-rate wireless systems. However, the focus is on developing mechanisms to ensure user satisfaction for the different data services for such systems from a network perspective. An initial fact finding of how the user satisfaction varies with heterogeneous wireless data services is conducted. Based on the findings we have focused on enhancing effective throughput, designing scheduling algorithms to ensure the strict time constraints specifically for multi-rate systems and implementing mobility management solutions for wireless LAN systems. To gain additional insight regarding the maximum possible number of satisfied users for a multi-rate wireless system, we have performed a theoretical study based on techniques derived from the auction theory. We conclude that individually the schemes do enhance the performance but only with joint functioning of the schemes that user satisfaction can be provided.

We first focused on capturing the variation of user satisfaction with varying network parameters. But prior to that we modeled the user satisfaction in order to transform a subjective quality into a quantitative metrics. We proposed user irritation metrics , both short and long term which reflected the user tolerance and sensitivity to the data service received. This study provides a deep understanding about the variation of user satisfaction with the network dynamics which actually effect the quality of service. Thereafter, we proposed a class-based QoS framework. It comprises of a radio resource management scheme which considers user satisfaction based on the perceived QoS, and caters to heterogeneous applications that have diverse QoS requirements. The proposed resource management scheme has two components: the admission control algorithm caters to the *long term* user satisfaction while the session-based rate and bandwidth allocation scheme

manipulates the *short term* user satisfaction.

The performance of the existing channel estimation schemes are limited due to the fact that they were based on the good and bad channel states and hence are incapable of harnessing the higher bandwidth provided by multi-rate wireless systems. In order to overcome the short comings, we developed a fast and accurate estimation technique based on information theoretic measures. The proposed technique though resulted in higher throughput but did not ensure that the strict timing requirements demanded by the real-time data services are ensured. Hence, as a next step, we developed scheduling algorithms based on the accurate channel estimation with the focus of assuring the strict timing requirements for the real-time data services in a multi-rate wireless system.

One of the inherent characteristics for wireless devices is mobility. Mobility solutions though well researched in the cellular domain is still a bottleneck for wireless lans. Consequently, the strict timing requirements for real-time wireless data services like VoIP, streaming multimedia etc. are not honored in current 802.11 deployments. In response to that we have designed, evaluated and implemented a new new mobility management framework that ensures *seamless* and *transparent* handoff between different access points in IEEE 802.11 based wireless local area networks. Our proposed solution is a client-end solution capable of reducing handoff delays in Wi-Fi networks to 15 ms at best, 20 ms in the average case and 26 ms in the worst case.

The concluding work of the dissertation is a theoretical study so as to determine the maximum possible number of satisfied users under a single access point for a multi-rate wireless system. Our proposed formulation focusses on guaranteeing the minimum data rate of the users, and as a secondary objective maximizes the overall system throughput. The analysis reveal that existing schemes suffer from the exposure problem and are not able to maximize the number of satisfied user. The proposed appraoch though theoer-tical can schedule more users whose minimum QoS requirements are met than existing

schemes. We have shown that the worst case performance of the proposed approximate algorithm is bounded by a multiplicative factor $(1 + \log m)$ corresponding to the optimal solution, where m denotes the number of slots in a schedule cycle. To conclude, we would like to emphasize that individually the proposed schemes do improve the performance of the multi-rate wireless systems, however to ensure the user satisfaction all the proposed schemes need to be employed and should work simultaneously.

6.0.1 Future Work

The dissertation has focussed only on the network aspects related to the multi-rate wireless systems. We feel that in order to ensure user satisfaction an indepth study of the impacts on the application layer is needed. Preliminary experiments have revealed that quality of multimedia streams suffer because the audio-video stream synchronization gets affected due to loss and delay variations at the lower layers.

In chapter 3, we have proposed estimation techniques and scheduling algorithms for multi-rate systems. Although we have been successful in improving performance however we do not have any insight on the length of the schedule vector, which would give the optimal performance. Simulation results though indicate that schedule length between 800 and 1000 seems to result in better system performance.

The mobility solutions presented in chapter 5 is limited to open authentication. Designing, implementing and upgrading the current solution to work in authentication mode would be challenging. The challenges will lie in authenticating between different networks during handoffs as well mitigating the extra delay that would be introduced due to the authentication operation.

REFERENCES

- [1] 3GPP TR 25.858 V5.0.0, “High Speed Downlink Packet Access: Physical Layer Aspects (Release5)”, Mar. 2002.
- [2] 3GPP2 C.S0024 Ver 3.0, “cdma2000 High Rate Packet Data Air Interface Spec.”, Dec. 5, 2001.
- [3] IEEE 802.11, Part II “Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications ”’, August 1999.
- [4] IS-856 cdma2000 High Rate Packet Data Air Interface Specification, *3GPP2 C.SO024 Version 4.0*, <http://www.3gpp2.org>, Oct 25, 2002.
- [5] “RTP payload format for H.261 video streams”, *Internet Engineering Task Force*, RFC 2032, Oct. 1996.
- [6] S. Aramvith, I. Pao and M. Sun, “A Rate-Control Scheme for Video Transport over Wireless Channels”, *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 11, no. 5, pp. 569-580., May 2001.
- [7] Atheros Communications , “<http://www.atheros.com>”.
- [8] L. Badia, M. Lindström, J. Zander, M. Zorzi, “Demand and Pricing Effects on the Radio Resource Allocation of Multimedia Communication Systems”, *Proceedings IEEE Globecom 2003*, San Francisco, CA, vol. 7, pp. 4116-4121, Dec. 2003.
- [9] G. Bao, “Performance evaluation of TCP/RLP protocol stack over CDMA wireless links”, *ACM Wireless Networks Journal*, vol. 2, 1996, pp. 229-237.
- [10] H. Balakrishnan, V. N. Padmanabhan, S. Seshan and R. H. Katz, “A Comparison of Mechanisms for Improving TCP Performance over Wireless Links”, *IEEE/ACM Transactions on Networking*, pp.756-769, vol. 5, Dec. 1997,
- [11] H. Balakrishnan and R. Katz, “Explicit Loss Notification and Wireless Web Performance”, *IEEE Globecom Internet Mini-Conference*, 1998.

- [12] S. Biaz, N. H. Vaidya, “Distinguishing congestion losses from wireless transmission losses: a negative result”, *International Conference on Computer Communications and Networks*, 1998, pp. 722-731.
- [13] P. Bender, P. Black, M. Grob, R. Padovani, N. Sindhushayana, and A. Viterbi, “CDMA/HDR: A Bandwidth-Efficient High-Speed Wireless Data Service for Nomadic Users”, *IEEE Communications Magazine*, July 2000, pp. 70-77.
- [14] M. Bender, S. Chakrabarti and S. Muthukrishnan, “Flow and stretch metrics for scheduling continuous job streams”, *Proc. ACM Symposium on Discrete Algorithms (SODA)*, pp. 270-279., 1998.
- [15] V. Bene, “General Stochastic Processes in the Theory of Queues”, *Addison-Wesley*, 1963.
- [16] V. Bharghavan, S. Lu, and T. Nandagopal, “Fair Queuing in Wireless Networks: Issues and Approaches”, *IEEE Personal Communications*, vol. 6, pp. 44-53, Feb. 1999.
- [17] T. Bonald and Alexandre Proutiere, “Wireless downlink data channels: user performance and cell dimensioning”, *Proceedings of the 9th annual international conference on Mobile computing and networking, (MobiCom)*, pp. 339-352, 2003.
- [18] S. Borst, P. Whiting, “Dynamic Rate Control Algorithms for HDR Throughput Optimization”, *Proceedings of IEEE INFOCOM*, vol. 2, pp. 976-985, 2001.
- [19] V. Brik, A. Mishra and S. Banerjee, “Eliminating handoff latencies in in 802.11 WLANs using Multiple Radios: Applications, Experience, and Evaluation”, *ACM IMC*, Oct. 2005.
- [20] S. Cen, P.C. Cosman, and G.M. Voelker, “End-to-end differentiation of congestion and wireless losses”, *Proc. Multimedia Computing and Networking Conference (MMCN)*, pp. 1-15, 2002.

- [21] H.K. Choi and J.O. Limb, "A Behavioral Model of Web Traffic", *International Conf. of Network Protocols (ICNP)*, pp. 327-334, 1999.
- [22] W. Chung, H. W. Lee, and J. Moon, "Downlink Capacity of CDMA/HDR", *Proc. IEEE 2001 Vehicular Technology Conference (VTC2001-Fall)*, pp. 1937-1941, vol. 3, 2001.
- [23] Darwin Streaming Server, <http://developer.apple.com/opensource/server/streaming/index.htm>
- [24] D. Erdogmus and J. C. Principe, "An Error-Entropy Minimization Algorithm for Supervised Training of Nonlinear adaptive systems", *IEEE Transactions on Signal Processing*, vol. 50, no. 7, pp. 1780-1786, July 2002.
- [25] D. Erdogmus and J. C. Principe, "Generalized Information Potential Criterion for Adaptive System Training", *IEEE Transactions on Neural Networks*, vol. 13, no. 5, pp. 1035-1044, 2002.
- [26] S. Ferrara, T. Matsumoto, M. Nicoli and, U. Spagnolini, "Soft Iterative Channel Estimation With Subspace and Rank Tracking", *IEEE Signal Processing Letters*, vol. 14, pp. 5-8, 2007.
- [27] M. Fine, K. McCloghrie, J. Seligson, K. Chan, S. Hahn, R. Sahita, A. Smith and F. Reichmeyer, *Framework Policy Information Base*, IETF-rap-frameworkpib-04.
- [28] F. Fitzek and M. Reisslein, "MPEG-4 and H.263 Video Traces for Network Performance Evaluation", *IEEE Network*, vol. 15, no. 6, pp. 40-54, 2001.
- [29] A. Furuskar, S. Mazur, F. Muller and H. Olofsson, "EDGE: Enhanced Data Rates for GSM and TDMA/136 Evolution", *IEEE Personal Communication Magazine*, pp 56-66, June 1999.
- [30] Gartner Research, "<http://www.gartner.com>".
- [31] General Packet Radio Services (GPRS) Service Description, *3GPP TS 03.60*, <http://www.3gpp.org>.

- [32] J. Hastad, "Clique Is Hard to Approximate within $n^{1-\epsilon}$ ", *Acta Mathematica*, vol. 182, pp. 105-142, 1999.
- [33] International Telecommunication Union, *G.114*, "One-way Transmission Time", May 2000.
- [34] H. M. Radha, M. Schaar and Y. Chen, "The MPEG-4 Fine-Grained Scalable Video Coding Method for Multimedia Streaming Over IP", *IEEE Transactions on Multimedia*, vol. 3, no. 1, pp. 53-68, 2001.
- [35] A. Jalali, R. Padovani, and R. Pankaj, "Data throughput of CDMA-HDR a high efficiency-high data rate personal communication wireless system", *IEEE Proceedings of Vehicular Technology Conference*, vol. 3, pp. 1854-1858, Spring 2000.
- [36] R. Koenen, "Overview of the MPEG-4 Standard", *ISO/IEC JTCLSC29/WG11 M4030*, 2001.
- [37] V. Krishna, "Auction Theory", *Academic Press*, USA, 2002.
- [38] S. S. Kulkarni and C. Rosenberg, "Opportunistic Scheduling Policies for Wireless Systems with Short Term Fairness Constraints", *Proceedings of IEEE Globecom*, vol. 1, pp. 533- 537, Dec 2003.
- [39] W. E. Leland, M. S. Taqqu, W. Willinger and D. V. Wilson, "On the self-similar nature of Ethernet traffic (extended version)", *IEEE/ACM Transactions of Networking*, vol. 2, no. 1, pp. 1063-6692, 1994.
- [40] X. Liu, E. K. P. Chong, and N. B. Shroff, "Transmission Scheduling for Efficient Wireless Utilization", *Proceedings of IEEE INFOCOM*, vol. 3, pp. 776-785, 2000.
- [41] X. Liu, E. K. P. Chong and N. B. Shroff, "Optimal Opportunistic Scheduling in Wireless Networks," *IEEE Vehicular Technology Conference (VTC Fall)*, Orlando, Florida, vol. 3, pp. 1417-1421, Oct. 2003.

- [42] X. Liu, E. K. P. Chong and N. B. Shroff, "Opportunistic Transmission Scheduling with Resource-Sharing Constraints in Wireless Networks", *IEEE Journal on Selected Areas in Communications*, vol. 19, no. 10, pp. 2053-2064, Oct 2001.
- [43] Y. Liu, S. Gruhl, and E. Knightly, "WCFQ: An opportunistic wireless scheduler with statistical fairness bounds", *IEEE Transactions on Wireless Communication*, vol. 2, Issue 5, pp. 1017-1028, Sept. 2003.
- [44] Y. Liu and E. Knightly, "Opportunistic Fair Scheduling over Multiple Wireless Channels", *Proceedings of IEEE INFOCOM*, vol. 2, pp. 1106-1115, 2003.
- [45] H. Lin, M. Chatterjee, S. K. Das and K. Basu, "ARC: An Integrated Admission and Rate Control Framework for CDMA Data Networks Based on Non-Cooperative Games", *Intl. Conference on Mobile Computing and Networking (MobiCom)*, 2003.
- [46] T. Ling and N. Shroff, "Scheduling Real-Time Traffic in Atm Networks", *Proceedings of IEEE INFOCOM* vol. 1, pp. 198-205, 1996.
- [47] S. Lu, V. Bharghavan, and R. Srikant, "Fair scheduling in wireless packet networks", *IEEE/ACM Transactions on Networking*, Aug. 1999, pp. 473-489.
- [48] MADWiFi Project, "<http://www.madwifi.org/>".
- [49] P. Maille, "Auctioning for Downlink Transmission Power in CDMA Cellular Systems", *Proceedings of 7th ACM International Symposium on Modeling, Analysis and Simulation of Wireless and Mobile Systems (MSWiM)*, pp. 293 - 296, Oct 2004.
- [50] P. Maille and B. Tuffin, "Multi-bid Auctions for Bandwidth Allocation in Communication Networks," *Proceedings of IEEE INFOCOM*, Mar 2004.
- [51] N.B. Mandayam, P.-C. Chen and J.M. Holtzman, "Minimum duration outage for cellular systems: a level crossing analysis", *Proc. of IEEE VTC*, pp. 879-883, 1996.
- [52] Y. Mansour, B. Patt-Shamir, "Jitter Control in QoS Networks", *IEEE/ACM Transactions on Networking*, vol. 9, no. 4, Aug. 2001, 1998.

- [53] V. Mhatre and K. Papagiannaki, "Using Smart Triggers for Improved User Performance in 802.11 Wireless Networks", *ACM Mobisys*, Jun. 2006.
- [54] Meeru Networks, "<http://www.merunetworks.com>".
- [55] Mplayer, www.mplayerhq.hu.
- [56] M. Nicoli, O. Simeone, and U. Spagnolini, "Multi-slot estimation of frequency-selective fast-varying channels", *IEEE Trans. Commun.*, vol. 51, no. 8, pp. 1337-1347, Aug. 2003.
- [57] Netgear WNIC, "<http://www.netgear.com/products/details/WAG511.php>".
- [58] S. Pack and Y. Choi, "Pre-Authenticated Fast Handoff in a Public Wireless LAN Based on IEEE 802.1x Model", *ACM PWC*, Oct. 2002.
- [59] J. Padhye, V. Firou, D. Towsley and J. Kurose, "Modeling TCP Throughput: A Simple Model and its Empirical Validation", *ACM SIGCOMM*, pp. 303-314, 1998.
- [60] S. Pal, M. Chatterjee and S. K. Das, "Improving Guarantees on Delivery Time for Wireless Data Services", *IEEE WCNC*, vol. 4, pp. 2539-2544, 2004.
- [61] S. Pal, M. Chatterjee and S. K. Das, "User-Satisfaction based Differentiated Services for Wireless Data Networks", *IEEE International Conference on Communications (ICC)*, May 2005.
- [62] S. Park, H. Kim, C. Park, J. Kim, and S. Ko, "Selective Channel Scanning for Fast Handoff in Wireless LAN Using Neighbor Graph", *In Proceedings of LNCS*, vol. 3260, January 2004.
- [63] E. Parzen, "On estimation of a probability density function and mode", *Time Series Analysis Papers*, 1967.
- [64] A. Pekec and M. H. Rothkopf, "Combinatorial Auction Design", *Management Science*, vol. 49, pp. 1485-1503, 2003.
- [65] I. Ramani and S. Savage. SyncScan: practical fast handoff for 802.11 infrastructure networks, *IEEE INFOCOM*, Mar. 2005.

- [66] Siegmund M. Redl, Matthias K. Weber, Malcolm W. Oliphant, "An Introduction to GSM", *Artech House*, Mar. 1995.
- [67] A. Renyi, *Probability Theory*, American Elsevier Publishing Company Inc., New York, 1970.
- [68] T. Sandholm, S. Suri, A. Gilpin and D. Levine, "Winner Determination in Combinatorial Auction Generalizations", *Proceedings of the First International Joint Conference on Autonomous Agents and Multiagent Systems*, pp. 69-76, Jul. 2002.
- [69] S. Shakkottai and A. Stolyar, "Scheduling Algorithms for a Mixture of Real-Time and Non-Real-Time Data in HDR", *Proc. International Teletraffic Congress (ITC-17)*, pp. 793-804, Sep. 2001.
- [70] M. Shin, A. Mishra and W. Arbaugh, "Improving the latency of 802.11 hand-offs using neighbor graphs", *In Proceedings of ACM/Usenix Mobisys*, Jun. 2004.
- [71] M. Shin, A. Mishra and W. Arbaugh. An empirical analysis of the IEEE 802.11 MAC layer handoff process, "*ACM SIGCOMM Comput. Commun. Rev.*", vol. 33, no. 2, 2003.
- [72] L. D. Soares and F. Pereira, "MPEG-4: a flexible coding standard for the emerging mobile multimedia applications", *IEEE International Symposium on Personal, Indoor and Mobile Radio Communications*, vol. 3, pp. 1335-1339, 1998.
- [73] Y. Snir, Y. Ramberg, J. Strassner, R. Cohen, *Policy Framework QoS Information Model*, IETF-policy-qos-info-model-02.
- [74] G.D. Stamoulis, D. Kalopsikakis and A. Kyrikoglou, "Efficient agent-based negotiation for telecommunications services", *Global Telecommunications Conference (GLOBECOM)*, vol 3. pp. 1989-1996, 1999.
- [75] J. Sun, L. Zheng and E. Modiano, "Wireless Channel Allocation Using an Auction Algorithm," *Proceedings of Allerton Conf. Communication, Control, and Computing*, vol. 24, pp. 1085-1096, Oct. 2003.

- [76] J. Sun, E. Modiano and L. Zheng, "A Novel Auction Algorithm for Fair Allocation of a Wireless Fading Channel", *Proceedings of 38th Annual Conference on Inform. Science and Systems*, Mar. 2004.
- [77] A. Tarello, E. Modiano, J. Sun, M. Zafer, "Minimum Energy Transmission Scheduling subject to Deadline Constraints", *Proceedings of 3rd IEEE Modeling and Optimization in Mobile, Ad Hoc, and Wireless Networks (WiOPT)*, pp. 67-76, 2005.
- [78] H. Velayos and G. Karlsson. Techniques to reduce the IEEE 802.11b handoff time, *IEEE ICC*, June 2004.
- [79] P. Vishwanath, D. N. C. Tse and V. Anantharam, "Opportunistic beamforming using dumb antennas", *IEEE Transactions on Information Theory*, vol. 48, Issue 6, pp. 1277-1294, Jun. 2002.
- [80] A. J. Viterbi, "Spread spectrum communications - myths and realities", *IEEE Communication Magazine*, May 1979.
- [81] N. Wiener, *Extrapolation, Interpolation, and Smoothing of Stationary Time Series with Engineering Applications*, MIT Press, 1949.
- [82] M. Xiao, N.B. Shroff, E.K.P Chong, "Utility-based power control in cellular wireless systems", *IEEE INFOCOM*, vol.1, pp. 412-421, Jun. 2002.
- [83] M. Yavuz, D. Paranchych, and G. Wu, "Performance Improvement of the HDR System due to Hybrid ARQ", *Proc. IEEE Vehicular Technology Conference (VTC2001-Fall)*, vol. 4, pp. 2192-2196, Oct. 2001.
- [84] W. Rae Young, "AMPS: Introduction, Background, and Objectives", *Bell System Technical Journal*, vol. 58, no. 1, pp. 1-14, Jan. 1979.
- [85] D. Zhao, X. Shen and J. W. Mark, "Radio Resource Management for Cellular CDMA Systems Supporting Heterogeneous Services", *IEEE Transactions on Mobile Computing*, vol. 2, no. 2, Jun. 2003.

- [86] M. Zorzi, "Outage and error events in bursty channels", *IEEE Transactions on Communications*, vol. 46, no. 3, pp. 349-356, Mar. 1998.

BIOGRAPHICAL STATEMENT

Sourav Pal received his Bachelors of Engineering in Computer Science and Engineering degree from Bengal Engineering College, Shibpor, India. He received his M.S. and Ph.D. degrees from The University of Texas at Arlington in 2003 and 2007, respectively, all in Computer Science and Engineering. His current research interests include computer networking with focus on wireless systems, multimedia systems and performance optimization. He has published more than a dozen research papers during his doctoral studies in these fields.