

# Dynamic QoS-Based Bandwidth Allocation Framework for Broadband Wireless Networks

Amir Esmailpour and Nidal Nasser, *Member, IEEE*

**Abstract**—Broadband wireless communication systems, namely, Worldwide Interoperability for Microwave Access (WiMAX) and Long-Term Evolution (LTE), promise to revolutionize the mobile users wireless experience by offering many of the services and features promised by fourth-generation (4G) wireless systems, such as supporting multimedia services with high data rates and wide coverage area, as well as all-Internet Protocol (IP) with security and quality-of-service (QoS) support. These systems, however, require proficient radio resource management (RRM) schemes to provide the aforementioned features they promise. In this paper, we propose a new framework, which is called dynamic QoS-based bandwidth allocation (DQBA), to support heterogeneous traffic with different QoS requirements in WiMAX networks. The DQBA framework operates as such; it dynamically changes the bandwidth allocation (BA) for ongoing and new arrival connections based on traffic characteristics and service demand. The DQBA aims at maximizing the system capacity by efficiently utilizing its resources and by being fair, practical, and in compliance with the IEEE 802.16 standard specifications. To achieve its objectives, DQBA employs a flexible architecture that combines the following related components: 1) a two-level packet scheduler scheme; 2) an efficient call admission control policy; and 3) a dynamic BA mechanism. Simulation results and comparisons with existing schemes show the effectiveness and strengths of the DQBA framework in delivering promising QoS and being fair to all classes of services in a WiMAX network.

**Index Terms**—Inter- and intraclass quality-of-service (QoS) support, packet scheduling (PS), QoS differentiation.

## I. INTRODUCTION

THE FUTURE of wireless networking and wireless communication systems will rely very much on the quality of service (QoS) they can provide to end users and the level at which they can satisfy application and service demands. Traditional high-capacity and long-range coverage for wireless applications will no longer dictate the success of such systems. New technologies such as Worldwide Interoperability for Microwave Access (WiMAX) and Third-Generation Partnership Project (3GPP) Long-Term Evolution (LTE) could provide the medium and facilitate the requirements for sophisticated applications such as multimedia, Internet Protocol television, remote surgery, and so on [1]. Radio resource management

(RRM) techniques such as resource allocation, scheduling, and admission control are among powerful strategies that could bring about such capabilities to future wireless technologies.

WiMAX is part of the next-generation broadband wireless access (BWA) technologies. BWA is emerging as an access network with several advantages. These include faster deployment, high scalability, low maintenance, as well as less cost and modular investment for upgrades. WiMAX is the most commonly used implementation of the IEEE 802.16 standard. It is highly attractive for its cost effectiveness and compatibility with fourth-generation all-Internet Protocol wireless networks.

WiMAX is designed to provide high-speed wireless access in metropolitan area networks. RRM techniques and QoS support are among the most important features of this technology. The IEEE 802.16 standard [2], which is associated with WiMAX, defines five classes of traffic flow representing different types of services in the following order: 1) unsolicited grant service (UGS); 2) extended real-time polling service (ertPS); 3) real-time polling service (rtPS); 4) nonreal-time polling service (nrtPS); and 5) best effort service (BE). The standard defines a connection-oriented medium-access control (MAC) protocol with a mechanism for QoS support. However, RRM techniques such as packet scheduling (PS), call admission control (CAC), and bandwidth allocation (BA) schemes are left open for development by the vendors [3].

Wireless technologies are facing challenging issues with respect to new applications and increasingly sophisticated user demands. Traditionally, the challenges for wireless systems involved sending large capacities of information across a long distance through wireless media. Most of the challenging issues dealt with wireless signal propagation problems such as noise, fading, interference, and so on. Although traditional problems such as the limited radio spectrum still remain among key problems, more sophisticated problems have evolved, which deal with QoS, resource allocation, and resource reservation.

Traffic generated by real-time (RT) and nonreal-time (NRT) applications with various demands and end users with diverse expectations are traversing a common and shared environment in the wireless network. Managing such complex systems and distributing resources among all users fairly and efficiently have proven to be a challenging task.

One key issue in providing multimedia services over a WiMAX network is providing QoS support to RT applications while avoiding starvation of NRT applications [4]. One way to mitigate this issue is to improve system utilization. However, another problem is to provide seamless service transition when a service class changes the QoS requirements dynamically. One solution to this challenge involves proper management of the

Manuscript received November 21, 2010; revised April 4, 2011; accepted May 6, 2011. Date of publication June 7, 2011; date of current version July 18, 2011. The review of this paper was coordinated by Dr. P. Lin.

A. Esmailpour is with the Department of Electrical and Computer Engineering, University of Toronto, Toronto, ON M5S 3G4, Canada (e-mail: Amir.Esmailpour@utoronto.ca).

N. Nasser is with the School of Computer Science, University of Guelph, Guelph, ON N1G 2W1, Canada (e-mail: nnasser@uoguelph.ca).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TVT.2011.2158674

resources to provide better QoS for end users and to achieve redistribution of system resources swiftly and smoothly. Designing an accurate and delicate QoS support management is the key to this solution.

Existing RRM techniques in the literature have several limitations in evaluating QoS requirements of the service classes, as well as the QoS deliveries to RT and NRT applications. In this paper, we address some of the limitations and propose respective solutions. One of the limitations with the existing solutions in PS and BA comes from separation of RT and NRT service requests and providing QoS to one group with no regard for the others. The IEEE 802.16 standard defines five categories of service flows [5]. Various RT and NRT applications fall into one of these categories. To be in compliance with the standard, WiMAX products and the corresponding RRM techniques ought to support both RT and NRT applications simultaneously.

One of the other challenging issues with the existing solutions comes from traffic separation. WiMAX allows the traffic flows to be separated in two levels, namely, intra- and interclass. In intraclass, the traffic is separated among several service flows within the same class of service. The QoS requirement of each flow is distinguished among other flows in the same class based on criteria set for that specific class of service. On the interclass level, traffic is separated among various classes of service. Service flows belonging to the same service class are grouped together, and their QoS requirement is presented as class based [6]. To have a clear distinction among each service flow from all service classes, an RRM strategy should distinguish the flows in both intra- and interclass levels simultaneously. Most of the research papers that we have seen dealt with service flows in either intra- or interclass levels but not both at the same time.

In this paper, we address some of the existing problems with providing QoS support for multimedia services and propose comprehensive solutions based on dynamic RRM techniques, including algorithms that control the usage of radio resources. To this end, we propose a framework for RRM that includes PS, CAC, and BA components and supports multimedia traffic with various QoS requirements such as priority, fairness, and utilization. The proposed framework for RRM in WiMAX is called "dynamic QoS-based BA" (DQBA). DQBA supports all types of service flows and makes a BA technique that is dynamic and fair and utilizes the resources in an efficient way. The motivation for this study is to develop RRM schemes that consider QoS differentiation and to establish parameters to quantify the differentiation based on fairness and utilization measurements.

The contributions of this paper include the following: 1) developing a QoS framework (including PS, CAC, and BA) that makes a delicate balance between the requirements of RT and NRT application demands. It improves the fairness for NRT applications while keeping the RT QoS guarantees at an acceptable level, 2) proposing parameters by which we can quantify QoS differentiation levels, and 3) incorporating QoS differentiation levels in the RRM decision-making process.

The components of the proposed DQBA framework are implemented in three overlapping building blocks called CAC, Tier-1, and Tier-2. The admission control policy is entirely implemented in the CAC module. PS is implemented in two stages, i.e., intraclass scheduling in Tier-1 and interclass in

Tier-2, and BA is mainly done in Tier-2. Each block includes several modules, and each module is responsible for achieving a set of goals. QoS differentiation, arrival rates, dynamic Cap (dCap), fairness, and utilization are all measured using their specific metrics. Each metric is calculated using different modules of Tier-1 and Tier-2 components. Inter- and intraclass traffic separation and QoS differentiation are also performed using two queuing modules in Tier-1 and Tier-2 components, respectively. The functional components of the framework are based on three fundamental ideas: 1) dynamically adjusting QoS support based on the knowledge of the network condition and the traffic behavior; 2) delivering QoS support in two levels, i.e., intraclass versus interclass; and 3) considering QoS differentiation while making the resource distribution decisions.

The rest of this paper is organized in the following order: Section II includes recent literature review in this area. Section III describes our proposed solution, including system architecture and modeling features. Section IV includes simulation and results. Finally, we make a conclusion and list future plans for this study in Section V.

## II. LITERATURE REVIEW

WiMAX has been recently labeled as one of the two contending technologies of the next generation of wireless networks. Among other special features of this technology, QoS and RRM have been in the forefront of the research activities. Numerous groups have proposed schemes for PS, CAC, and BA strategies. RRM schemes are categorized based on centralized versus decentralized, RT versus NRT, uplink versus DL, or basic versus complex schemes [7].

In the basic schemes, often a simple traditional queuing discipline is considered for PS of service classes [8]. Such systems use one or two simple queuing disciplines to separate the traffic. To improve the performance of these schemes, giving different weights to various classes to highlight their priorities is used as a common strategy. Wongthavarawat and Ganz [8] proposed two different PS with BA schemes for the IEEE 802.16 standard based on the weighted fair scheduling and the throughput guarantee scheduling schemes.

Proposals based on complex schemes for RRM schemes combine several queuing disciplines with some parameters related to traffic and medium characteristics. They often develop schemes in more than one stage. Complex hierarchical schemes were proposed in [9] and [10]. In hierarchical schemes, PS and BA are designed in multiple levels. Often, the traffic from different service classes is separated in the first level and then scheduled within each class in the second level. Different components of RRM could all be performed in the base station (BS), or the tasks could be divided between the BS and subscriber stations (SSs). In [11], the authors proposed an ad hoc scheme for scheduling, where a separate queuing discipline is used for each service flow.

Liang *et al.* proposed a novel scheduling scheme in [12], which covers both RT and NRT applications that perform resource allocation based on a polling interval adjustment. They show by simulation that their scheme could improve the performance of nrtPS while keeping that of rtPS in a steady state.

Bai *et al.* proposed a new QoS control protocol [13] for a point-to-multipoint (PMP) mode of operation. Their scheme enables control of the QoS guarantee per connection, and by distributing the tasks to the SSs, they manage to reduce the signaling overhead. Their QoS scheme is designed based on a cross-layer method, which makes it robust against wireless link degradation. Their simulation results reveal that the QoS requirements of rtPS and nrtPS connections are maintained for the minimum reserved and the maximum sustained traffic rates. This scheme provides good support for polling services; however, for other RT service classes, it shows that the requested service levels drop. Although the authors claim that the QoS thresholds are not jeopardized, their results show that the response to the bandwidth starvation of connection also drops during link degradation. These results reveal that this scheme does not provide a satisfactory level of fairness across various service classes.

Belghith proposed a scheduling scheme in [14] for the QoS classes in WiMAX, which considers the pricing factor in the optimization of resource usage and resource allocation. The optimization of RRM provides higher revenue for the service providers by showing the effects of the pricing model on the scheduling scheme. This scheme is a good indication of how utilization could impact the overall performance and, thereby, the revenue associated with the operation. The author proposed different pricing schemes for various classes of service, which also indicates that fairness in allocation of resources could eventually affect the overall revenue and justification of economic aspects of the operation.

Pizzi *et al.* used deficit round-robin (DRR)-based schedulers in [15] for the MAC layer of a PMP mode of operation in the WiMAX network and compared two methods of a compensation-based algorithm versus a greedy algorithm. Their simulation results show that the compensation-based approach is more capable of service differentiation when it comes to heterogeneous channel conditions. In [16], Pizzi *et al.* further investigated the effect of adaptive modulation and coding schemes on the performance of the scheduling algorithms and argued that the evaluation of such algorithms is based on their ability to provide QoS differentiation among traffic classes and fairness in the treatment of data flows in the same classes. Although adaptive modulation coding is used to dynamically adapt to channel conditions, it does not guarantee interclass QoS differentiation and intraclass fairness. Queues are managed by the scheduler on a per-frame basis using a simple DRR scheduling algorithm for polling services while serving UGS and ertPS based on constant resource reservation. Their design is based on evaluation of opportunistic versus compensation-based techniques, and they conclude that this approach does not provide an effective strategy for service differentiation and fairness among various service flows. Their results show that channel awareness is mandatory in WiMAX for both modulation and coding schemes and to improve the scheduler performance. This paper does not provide a clear explanation with respect to the fact that the opportunistic methods could yield to unfairness caused by starvation of traffic flows coming with poor channel conditions, whereas compensation-based methods are capable of distributing the available bandwidth with fairness.

The authors in [17] studied the behavior of different scheduling algorithms, with main focus on rtPS applications. They also proposed a new scheduling scheme that could provide them with less mean delay compared with other proposed schemes such as the temporary removal scheduler (TRS) [18] and the maximum signal-to-interference ratio (mSIR) [19]-based scheduler. Their proposed scheduler serves the SSs with a minimum SNR after serving SSs with a greater SNR, resulting in less delay and higher throughput. The scheme also gives priority to SSs with a higher signal-to-interference ratio. They show by simulation that their proposed scheme could provide enhancements to throughput and mean delay, as well as delivered data packets per frame, as compared with both TRS and mSIR. This study is only based on an rtPS service class, and like other SNR-based schedulers, it cannot guarantee fairness to a distanced subscriber and does not provide improvements in this area, as compared with other related studies.

From a different perspective, RRM could be classified based on QoS differentiation strategies. Most of the proposed solutions that we have seen in this area are divided into two major groups, i.e., fairness- or priority-based models. Examples of these proposals are based on modified DRR (MDRR) [20], [21] and modified priority queue (MPQ) [11], [22] for fairness-based versus priority-based solutions, respectively. In this paper, we propose DQBA and compare the results of performance evaluation from these models with those of our proposed solution. DQBA uses several strategies such as multilevel queuing discipline, QoS differentiation, and unused bandwidth reallocation to improve the network performance.

The QoS differentiation in this paper is based on two important performance metrics, namely, fairness and utilization. Fairness in communication networks is used to evaluate whether the users or applications are receiving their fair share of resources. In the case of WiMAX applications, we refer to fairness in relation to BA among different service classes. There are several definitions and methods for evaluating fairness, such as Jain's fairness, the Gini coefficient, and max-min fairness [19]. Utilization in the wireless networks is defined as the ratio of the used network resources to the available resources. In next-generation heterogeneous multiservice wireless systems, resource utilization depends on several factors, including application demands, user density, and available network resources. In [23], utilization in terms of bandwidth consumption is defined. In this paper, we define fairness and utilization as the ratio of the allocated bandwidth by the requested bandwidth and the throughput achieved by the allocated bandwidth, respectively.

Chuck and Chang proposed in [24] the bandwidth recycling idea, which recycles the unused bandwidth while keeping the existing bandwidth reservation intact. This is similar to the bandwidth reallocation component of DQBA, which allows other SSs to utilize the unused bandwidth reclaimed from other service requests. This scheme provides QoS guarantees while improving system utilization and throughput. Their simulation results show that they can recycle 20% of the unused bandwidth in a first stage. Furthermore, they propose three scheduling algorithms to subsequently improve the throughput in two more stages to 26% and 30% when the network is in the steady state. Although the authors claim that they also manage to reduce



the delay with negligible overhead and without degrading the QoS requirements, the three-stage recycling algorithm seems to introduce new complexity to the already lengthy scheduling algorithms. Such complexity could very well introduce large overhead, thus subsequently degrading the QoS guarantees. The authors did not address the overheads generated and possible QoS degradation caused by that.

### III. DYNAMIC QUALITY-OF-SERVICE-BASED BANDWIDTH ALLOCATION SYSTEM DESIGN

In this paper, we propose a framework for the three components of RRM in WiMAX, which provides QoS with support for different types of applications supported by this technology. The proposed DQBA makes the BA for the time-division duplex mode of operation of the IEEE 802.16 standard and supports all types of service flows as indicated by the standard. DQBA is fair, efficient, and granular, and it follows the IEEE 802.16 variable transmission rate and signaling mechanism rules and regulations through dynamic adjustments of the BA.

DQBA is implemented in three building blocks to perform PS, CAC, and BA for WiMAX classes of service: CAC, Tier-1, and Tier-2. Each block includes several modules performing the tasks for scheduling, admission control, and BA. Details of each block and its modules are shown in Fig. 1.

#### A. PS Scheme

In Tier-1, four different algorithms are employed for five classes of service to perform intraclass scheduling, as shown in Fig. 2. For all connections holding UGS traffic, the scheduler uses earliest deadline first (EDF), where packets with the stringent deadlines will be scheduled first. The scheduler determines the packet's deadline based on its arrival time and maximum latency. EDF uses a dynamic priority-based scheduling algorithm that provides delay-guaranteed service to the UGS service class.

For ertPS and rtPS connections, weighted fair queuing (WFQ) is employed, which is similar to the generic fair queuing algorithm. The advantage of WFQ is that it is based on the resource reservation, i.e., instead of giving an equal share of resources to all users, it allocates a specific amount of resource to each service. However, it does not impose a strict time limit on each packet. ertPS and rtPS packets are scheduled based on their weights, which are calculated using the ratio between a connection's reserved traffic rate and the total sum of all reserved traffic rates of all connections.

For nrtPS connections, the scheduler uses round-robin (RR) queuing discipline, where each service gets a fair share of the allocated bandwidth in an RR fashion. RR assigns time slots to each service flow in equal portions and in order, handling all requests without priority. RR is both simple and easy to implement, and it does not cause any starvation problem for the lower priority services. BE connections will be treated in a first-in first-out (FIFO) fashion. FIFO is the most basic discipline, yet it is suitable for BE traffic with no particular QoS requests.

The scheduler goes through all queues in the order of their priorities, i.e., from the highest priority access to the lowest, and checks for any activities. If there is a packet waiting in a queue,

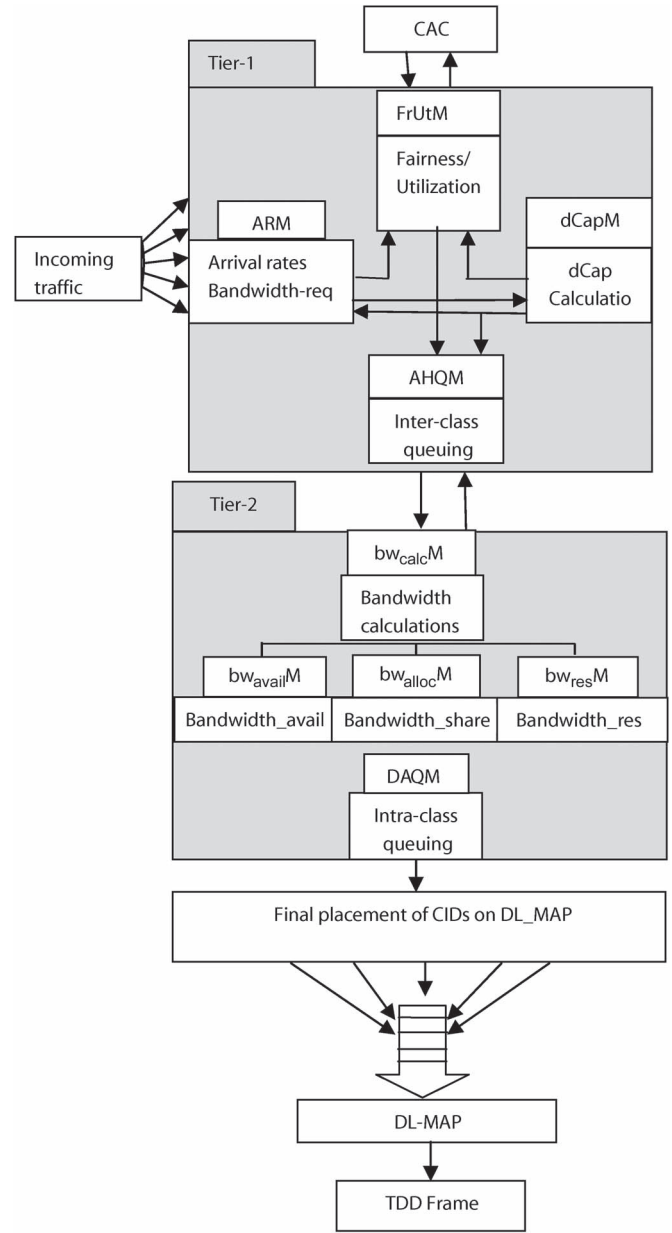


Fig. 1. DQBA framework architecture.

the queue will be added to an active queue list (AQL). If a queue belongs to the AQL in Tier-1, then the scheduler carries on to the final stage of scheduling and forwards packets to Tier-2 or the dynamic allocating queue (DAQ) stage to perform interclass scheduling. DAQ will dynamically assign a portion of the data burst for each type of traffic. If a queue has no packets, it will be removed from the AQL, and no bandwidth is allocated to it.

The process goes through several iterations. In each iteration, it updates the AQL, the available frame portion, and the number of packets in each queue. Then, it goes through several steps for each service round, processing higher priorities first and then the lower priorities. In each service round, the highest priority queue (PQ) will be served first until the available bandwidth is lower than the requested bandwidth. If the available bandwidth is lower than the bandwidth requested by the higher PQ, then the lower PQs will be served in order, providing a fair solution for the lower PQs.

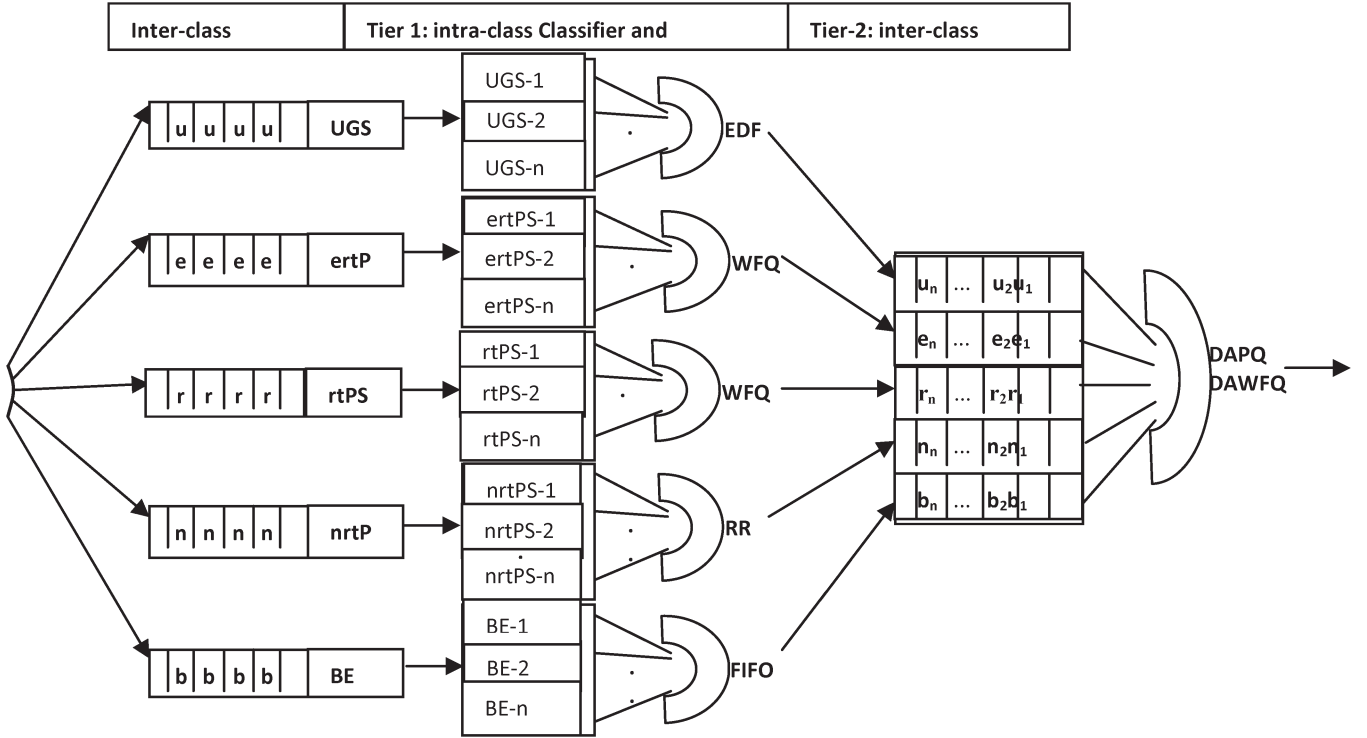


Fig. 2. 2-Tier scheduling scheme system design.

Tier-2 or DAQ has two types of implantation, i.e., dynamically allocating PQ (DAPQ) and dynamically allocating weighted fair queue (DAWFQ). In the case of DAPQ, traffic flows are separated based on their priorities, whereas in DAWFQ, they are set apart on a fairness basis. However, in both cases, QoS is differentiated using relative dCap (defined in the next section) values either as a priority factor for DAPQ or as weights of DAWFQ, respectively.

DAPQ is a modified version of the PQ. The problem with the PQ is that the highest priority will get allocated bandwidth, and if the high-priority flow continues for a long time, the lower priorities such as BE suffers from bandwidth starvation. We propose DAPQ as a solution for this problem and a fairness strategy based on a reservation technique.

DAWFQ is based on modified version of WFQ to allocate resources to different classes of service. DAWFQ schedules the packets based on the weight of each service class, which is calculated using the corresponding dCap and fairness values for each service class. In this case, we use interclass fairness. We dynamically calculate the weights based on fairness to avoid oversubscription of higher priority classes. In each round of scheduling, class-level fairness is considered, as is the total fairness.

### B. BA Strategy

The proposed dynamic BA (DBA) adjusts resource allocation based on traffic behavior and network conditions. The traffic behavior is characterized by the traffic arrival rate, whereas the network conditions are incorporated using two performance metrics, namely, fairness (Fr) and utilization (Ut).

The traffic arrival rate is used to calculate a new parameter called dCap. dCap is an important parameter used for PS, CAC, and DBA by integrating traffic and network characteristics and allows the system to improve fairness and utilization. dCap is a delicately calculated parameter that takes into account traffic behaviors such as traffic arrival rate, as well as network conditions such as fairness and utilization. The dCap value is calculated in (1). Calculation of effective rate  $\lambda_{\text{eff}}$ , fairness Fr, and utilization Ut factors, as stated in (1), is performed in various modules of the framework. A detailed explanation of the functionality of the modules in Tier-1 and Tier-2 are clearly outlined in Fig. 1. Thus

$$A \equiv \left( \text{Fr}^{n,j} < \text{Fr}_{\text{rel}}^{n,j-1} \right)$$

$$B \equiv \left( \text{Ut}^{n,j} < \text{Ut}_{\text{rel}}^{n,j-1} \right)$$

$$\text{dCap}^{n,j} = \begin{cases} \lambda_{\text{eff},i}^{n,j}, & \text{if } A \text{ and } B \\ & \text{are not true} \\ \lambda_{\text{eff},i}^{n,j} \left( 1 + \text{Fr}_{\text{rel}}^{n,j} \right), & \text{if } A \text{ is true} \\ \lambda_{\text{eff},i}^{n,j} \left( 1 + \text{Ut}_{\text{rel}}^{n,j} \right), & \text{if } B \text{ is true} \\ \lambda_{\text{eff},i}^{n,j} \left( 1 + \text{Fr}_{\text{rel}}^{n,j} + \text{Ut}_{\text{rel}}^{n,j} \right), & \text{if } A \text{ and } B \\ & \text{are true} \end{cases} \quad (1)$$

where  $\lambda_{\text{eff},i}^{n,j}$  is the normalized effective traffic rate for the connection  $i$  in the  $j$ th round of scheduling for service class  $n$ , and Fr and Ut are the corresponding fairness and utilization factors. dCap is used for both DAPQ and DAWFQ implementations of the Tier-2 component of the proposed solution, and it allows us to dynamically control the amount of resources allocated

based on the most recent network and traffic condition. The pseudocode for calculation of dCap is presented in Algorithm 1.

**Algorithm 1: dCapM algorithm for dCap calculation**

```

1: Check incoming packets
2: GET effective rate for the service request //ARM module
   (see Fig. 1)
3: GET fairness and utilization for the service request
   //FrUtM module
4: GET fairness and utilization for the corresponding class
5: For each service request
6: if ( $Fr^{n,j} \geq Fr_{rel}^{n,j}$ ) and ( $Ut^{n,j} \geq Ut_{rel}^{n,j}$ )
7:   then ( $dCap^{n,j+1} = \lambda_{eff}^{n,j}$ )
8: elseif ( $Fr^{n,j} < Fr_{rel}^{n,j}$ )
9:   then ( $dCap^{n,j+1} = \lambda_{eff}^{n,j} (1 + Fr_{rel}^{n,j})$ )
10: elseif ( $Ut^{n,j} < Ut_{rel}^{n,j}$ )
11:   then ( $dCap^{n,j+1} = \lambda_{eff}^{n,j} (1 + Ut_{rel}^{n,j})$ )
12: elseif ( $Fr^{n,j} < Fr_{rel}^{n,j}$ ) and ( $Ut^{n,j} < Ut_{rel}^{n,j}$ )
13:   then ( $dCap^{n,j+1} = \lambda_{eff}^{n,j} (1 + Fr_{rel}^{n,j} + Ut_{rel}^{n,j})$ )
14: Endfor
15: Return dCap value

```

In this paper, we define fairness as the ratio of the allocated bandwidth divided by the requested bandwidth for each service flow, and by summing up this parameter over a period of time for all flows within a service class, we get the total fairness of the system toward a service class  $n$  according to

$$Fr_T^{n,j} = \sum_{i=1}^I \left( \frac{bw_{alloc,i}^{n,j}}{bw_{req,i}^{n,j}} \right) \quad (2)$$

where  $i$  is the connection ID for the flow id  $i$  arriving in the  $j$ th round of scheduling for service class  $n$ .

Utilization, on the other hand, is the ratio of the throughput achieved divided by the allocated bandwidth, and it is calculated for all service flows in class  $n$  based on

$$Ut_T^{n,j} = \sum_{i=1}^I \left( \frac{Th_i^{n,j}}{bw_{alloc,i}^{n,j}} \right). \quad (3)$$

Upon calculation of dCap and corresponding Fr and Ut values in Tier-1, then the BA for the next round of scheduling is performed in the Tier-2 block according to

$$bw_{alloc}^T = \left( \frac{dCap^n}{dCap^T} \right) \times bw_{avail}^T \quad (4)$$

$$bw_{avail}^T = \sum_{n \in N} \sum_{i=1}^I bw_{res,i}^{j-1} + bw_{rem}^T \quad (5)$$

$$bw_{res,i}^{n,j} = bw_{alloc,i}^{n,j-1} - bw_{req,i}^{n,j} \quad (6)$$

$$bw_{rem}^T = BW - bw^T \quad (7)$$

where

$BW$  total link capacity;  
 $bw^T$  total occupied bandwidth;  
 $bw_{rem}$  remaining bandwidth;  
 $bw_{res}$  residual bandwidth;

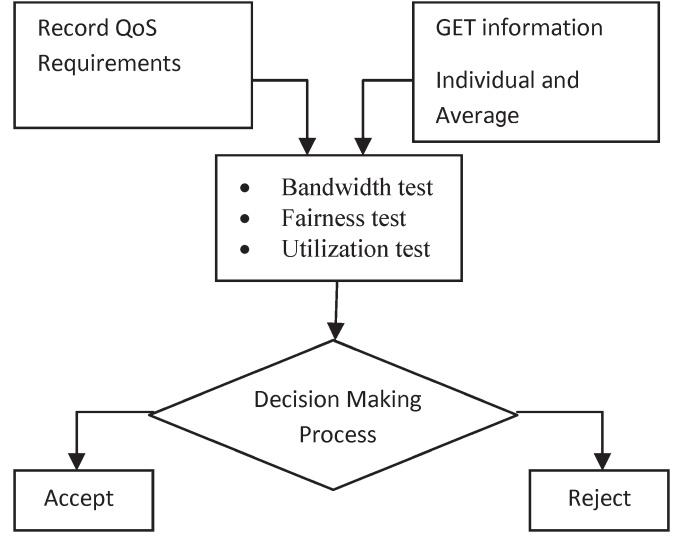


Fig. 3. CAC policy flowchart.

$bw_{req}$  bandwidth request;  
 $bw_{alloc}$  bandwidth allocated;  
 $bw_{avail}$  available bandwidth.

Allocated bandwidth  $bw_{alloc}^T$  is assigned to each service flow based on the proportion of available bandwidth  $bw_{avail}^T$  to the service request. This ratio is calculated based on the dynamic relative value of dCap ( $dCap^n/dCap^T$ ).  $bw_{avail}^T$ , in turn, is calculated based on the total of remaining bandwidth  $bw_{rem}^T$  plus the sum of total residual bandwidth  $bw_{res,i}^{j-1}$  in the previous round of scheduling. The residual bandwidth is the portion of the bandwidth allocated to a service class that is not being used, and the remaining bandwidth is the total capacity of the link minus the portion of the total bandwidth that is currently occupied.

### C. CAC Policy

CAC architecture includes a main structure that accepts the connection requests from each service flow, checks the conditions based on other components, and makes the final decision to whether accept or reject the call. CAC performs eligibility of new calls based on resource request, resource availability, and traffic behavior, as well as network conditions such as traffic arrival rate and fairness, respectively. An illustration of the overall process flow of CAC module is shown in Fig. 3.

CAC checks the bandwidth request and the bandwidth available for each service flow. If the bandwidth test has passed, then it performs two other tests, i.e., fairness and utilization tests. If all of the three tests have passed, then the call is accepted; otherwise, the call is rejected. Fairness and utilization tests are based on the history of fairness and utilization for this service class in the last round of scheduling.

## IV. SIMULATION AND RESULTS

We simulated a simple WiMAX network in OPNET Modeler version 14.5-PL3. The network consists of a single cell-based structure containing one BS and five SSs, as shown in Fig. 4.

The source of traffic is an application server that provides five types of applications, one for each type of traffic corresponding to five classes of service. The assumption is that each SS carries

TABLE I  
APPLICATIONS AND TYPES OF TRAFFIC AND QoS REQUIREMENTS

Service class	Class 1	Class 2	Class 3	Class 4	Class 5
Traffic type	UGS	ertPS	rtPS	nrtPS	BE
Application	VoIP	Video	SSH	FTP	browsing
Priority level	5 (highest)	4	3	2	1 (lowest)
Bandwidth requirement	5.0 Mbps	3.0 Mbps	1.5 Mbps	64 Kbps	32 Kbps
Delay Tolerance	10 ms	50 ms	50 ms	200 ms	500 ms

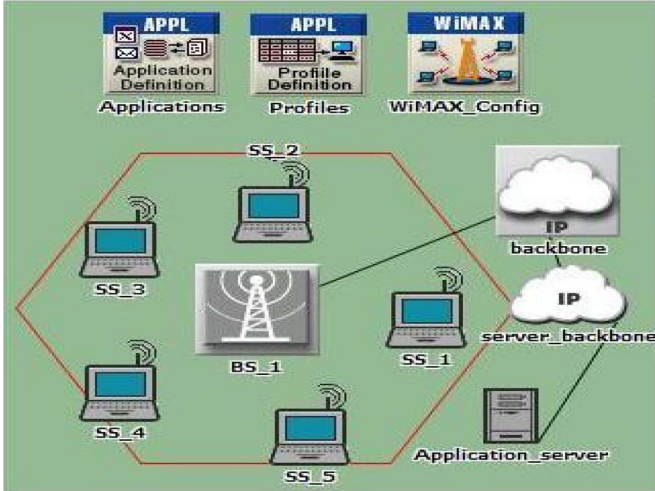


Fig. 4. OPNET simulation of a small WiMAX network including one cell, one BS, and five SSs.

the traffic from only one user, and each user is using only one type of application at a time. Applications, the type of traffic they represent, and subscribers that request them are listed in Table I. The values for bandwidth requirements and delays presented in Table I are typical industry benchmarks for these types of applications [25], [26].

#### A. Performance Comparison of DQBA With Other Models

The simulation starts by downloading files from the application server, one by each SS with the DQBA scheme in place. In the first case scenario, we look at the traffic separation and resource allocation for all five types of traffic using the proposed solution. The results of the throughput achieved by each type of service are shown in Fig. 5. The results in this figure show that DQBA has separated the traffic from all sources and delivered QoS requirements of all service requests accordingly. Shortly after the beginning of the simulation, the UGS throughput reaches its fixed bandwidth allocated at around slightly below 5 Mb/s, and subsequently, other RT classes reach their respective BAs, with ertPS at around 2.5 Mb/s and rtPS at 1.1 Mb/s. As for NRT traffic, nrtPS and BE have achieved almost 64 kb/s and over 30 kb/s, respectively. The results indicate that DQBA delivered resource request to NRT service flows fairly while maintaining the QoS requirements of RT applications in terms of bandwidth requirements.

The results show that DQBA satisfies the QoS requirement of RT classes while not overlooking the NRT classes. The results in Fig. 6 show that the delay value for UGS is less than 10 ms, which falls into the acceptable delay tolerance for the UGS class. For other RT applications, the delays have reached levels

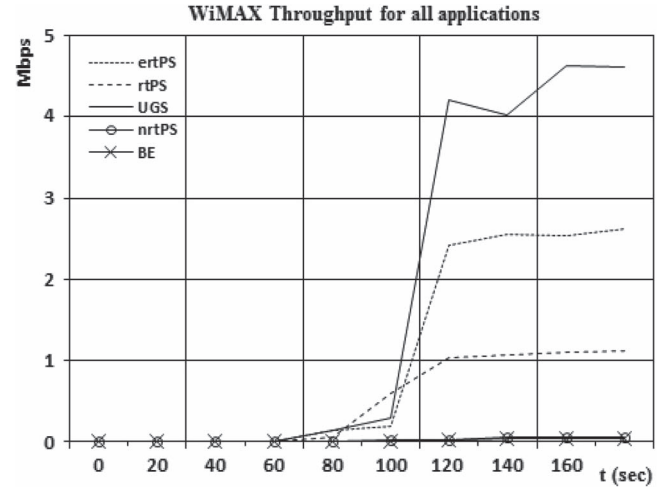


Fig. 5. Throughput achieved by all applications using DQBA.

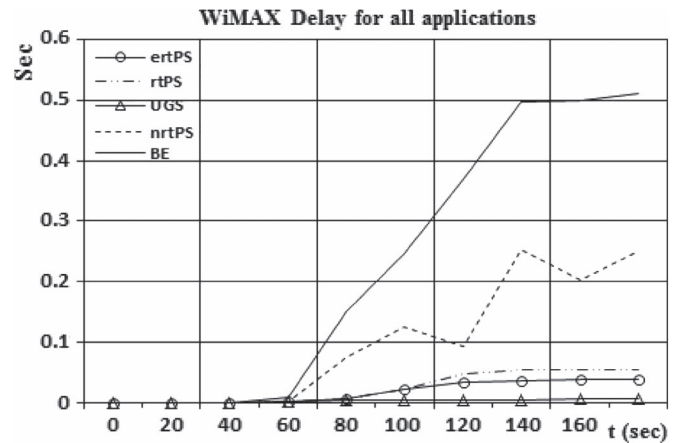


Fig. 6. Delay achieved by all applications using DQBA.

above 30 and 50 ms for ertPS and rtPS flows, respectively, for the duration of the simulation. The results for RT applications are slightly above the expected values. This is contributed to additional queuing delays generated by the NRT application demands. The delay results for NRT applications are also above 200 and 500 ms for nrtPS and BE traffic, respectively. However, for NRT applications, delay is not as critical as BA, which is quite satisfactory.

These results indicate that although the low-priority flows are not completely starved, the higher priority RT applications are consuming a considerable portion of the bandwidth. In this scenario, a link with a high data rate is employed, resembling a typical WiMAX overprovisioning case, in which the total bandwidth is not completely utilized.

In the second scenario, we repeat the same experiments with MPQ [27]. MPQ is a modified version of PQ, which considers



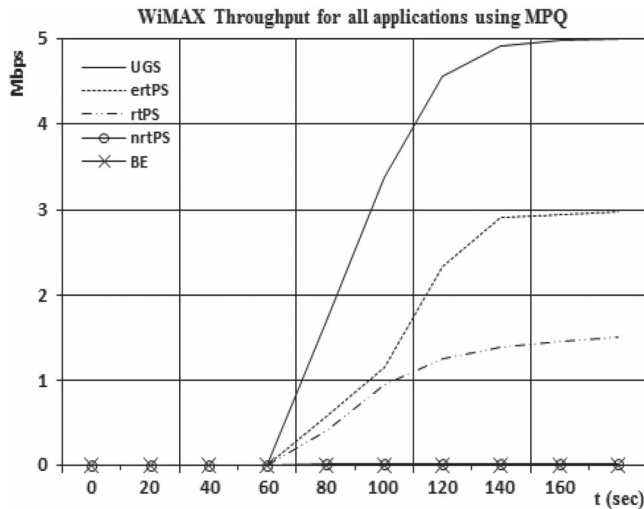


Fig. 7. Throughput achieved by all applications using MPQ.

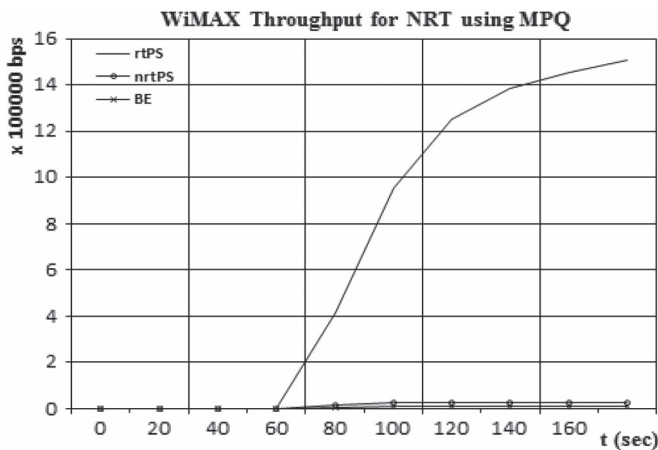


Fig. 8. Throughput for NRT and rtPS classes using MPQ.

the higher priority service classes with a modification that controls the amount of BA based on a maximum allowable value. Although the maximum allowable values control the amount of bandwidth, the service classes from the RT applications still take precedence over service flows from other classes. The results are shown in Figs. 7 and 8.

Fig. 7 shows that MPQ provides throughput for UGS and RT applications at levels close to bandwidth requirements and at a faster pace. UGS reaches 5 Mb/s quickly after the start of the simulation period. Other RT applications follow the same trend at 3 and 1.5 Mb/s for ertPS and rtPS, respectively.

To investigate further, we look at the throughput results for the NRT traffic flows using MPQ in Fig. 8 and compare them with those of DQBA. Fig. 8 shows the throughput values for NRT service flows in comparison with the lowest priority RT application (i.e., rtPS). As shown in this figure, the average RT application throughput is in the megabit-per-second ranges, which is close to 1.5 Mb/s, whereas the average NRT applications are in the lower ranges of kilobits per second, which do not match the bandwidth requested by these service classes.

When comparing these results with those of DQBA, it is observed that, in terms of RT applications, both models show high values in the megabit-per-second ranges; however, in terms of

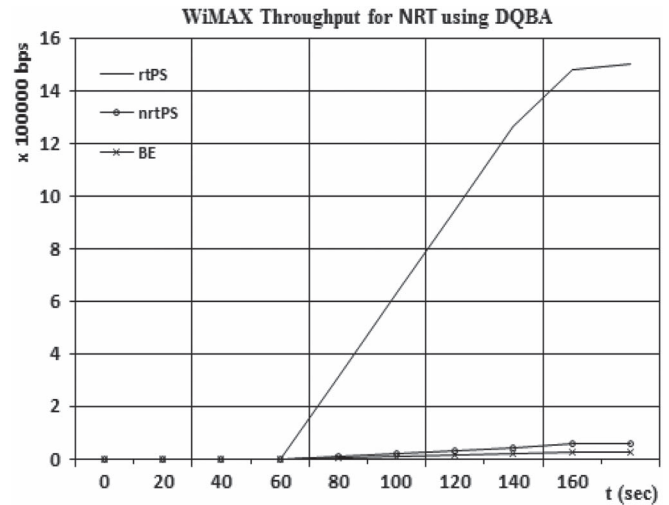


Fig. 9. Throughput for NRT and rtPS classes using DQBA.

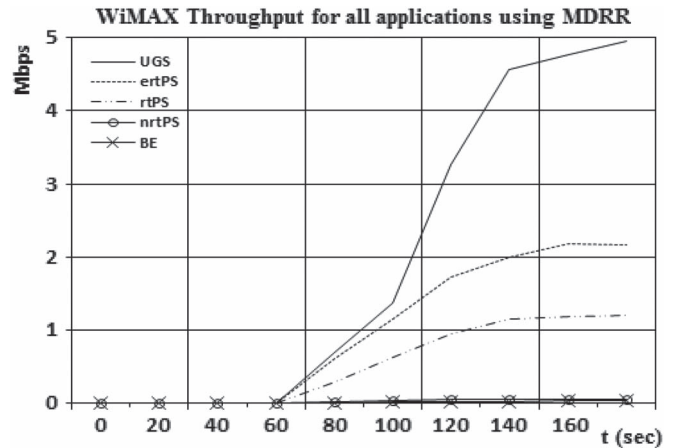


Fig. 10. Throughput for all applications using MDRR.

NRT applications, DQBA shows significant improvement in Fig. 9. Values for nrtPS and BE using DQBA in Fig. 9 are at 63 and 30 kb/s, as compared with those of MPQ in lower ranges of 17 and 10 kb/s, respectively, in Fig. 8.

According to the results from Figs. 5–9, RT traffic flows reach their respective BAs in MPQ faster than the proposed DQBA scheme. The BAs and throughputs are in comparable ranges in both cases. However, in terms of NRT applications, the proposed solution proves superior in comparison with MPQ while keeping the RT QoS requirements at a satisfactory level.

Fig. 10 shows a similar comparison among all service classes using the MDRR method. MDRR is a modified version of DRR, which takes into consideration a weight factor in addition to the deficit counter of DRR. The weight could be assigned a value based on the packet size of specific applications or a factor related to traffic, such as the bandwidth requirement or the packet arrival rate. In this case, we have associated the bandwidth request of each service class to the deficit counter as weights for this service class.

Based on the results in this figure, MDRR can provide the fixed bandwidth requirement for UGS while reducing the throughput for other RT applications (ertPS and rtPS) to above 2 and 1 Mb/s, respectively. Fig. 11, on the other hand, shows that



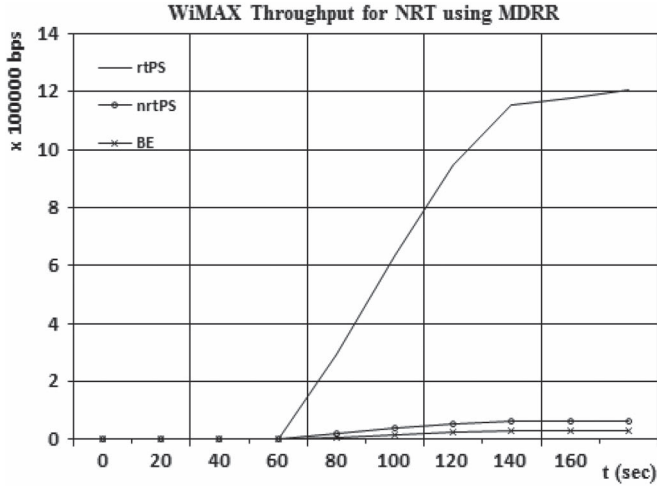


Fig. 11. Throughput for NRT classes using MDRR.

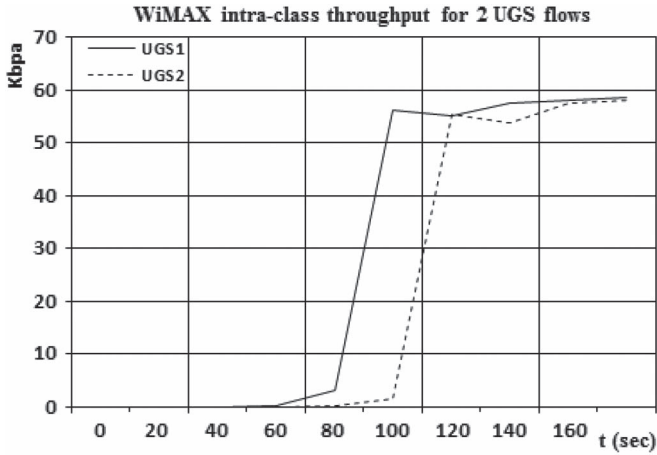


Fig. 12. Throughput for two UGS service flows (intraclass).

the resource allocations to NRT applications (nrtPS and BE) are significantly improved to around 60 kb/s and just below 30 kb/s when using MDRR in comparison to MPQ. These results were expected due to the nature of MDRR, by favoring NRT applications and improving the QoS guarantees to those service classes. However, the clear cost of such improvements is the reduction in throughput values for RT applications, as shown in Fig. 10. This figure shows that DQBA, in comparison with MDRR, provides elevated QoS support in terms of RT application demands while keeping a compatible level of support for NRT.

The overall results show that the proposed solution has a superior performance in comparison with MPQ when providing QoS support to NRT applications while keeping a steady level of support to RT applications. In addition, the proposed solution outperforms the MDRR solution in terms of RT applications with reasonable support for NRT applications.

### B. Intraclass Versus Interclass QoS Support

DQBA separates the traffic into two levels, i.e., intra- and interclass levels, and guarantees the QoS support at both levels using Tier-1 and Tier-2 components. Fig. 12 shows that the

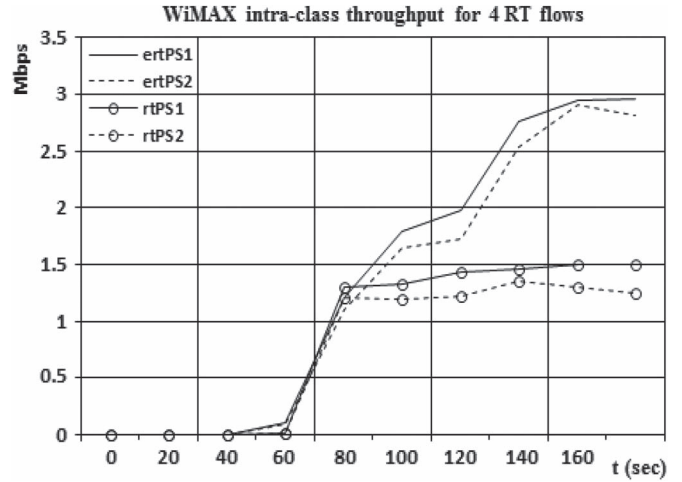


Fig. 13. Throughput comparison for two ertPS and two rtPS flows separated intraclass.

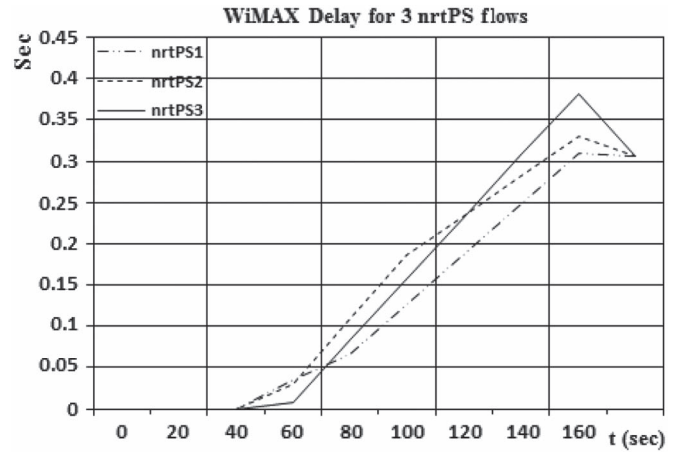


Fig. 14. Delay for three nrtPS flows separated intraclass.

traffic is separated intraclass between two UGS traffic flows with different delay tolerance values, using the EDF queuing discipline in the Tier-1 component of the proposed framework. Fig. 12 also shows that different QoS requirements for both UGS1 and UGS2 flows are met, however, with different delays. The throughputs achieved show that the bandwidth requirements are met at close to 5 Mb/s values, as shown in Table I.

In Fig. 13, two flows from ertPS classes along with two flows from rtPS classes with different weights (based on their relative dCap values) are separated according to their weights using the WFQ queuing discipline in Tier-1. Throughput values for ertPS1 and ertPS2 are in ranges close to 3 Mb/s, which indicate that the acceptable bandwidth requirements for an ertPS service class have been provided. The weights associated with flows ertPS1 and ertPS2 are calculated to be 0.23 and 0.24 for ertPS1 and ertPS2, and 0.18 and 0.19 for rtPS1 and rtPS2, respectively. Although a clear correlation between the weights and QoS requirements requires further studies beyond the scope of this paper, Fig. 13 shows that the proposed framework satisfies the intraclass QoS requirements by allocating various levels of resources, which are calculated by corresponding weights.

Fig. 14 shows the delay results for three NRT flows separated in an intraclass level using DQBA. NRT flows are separated,

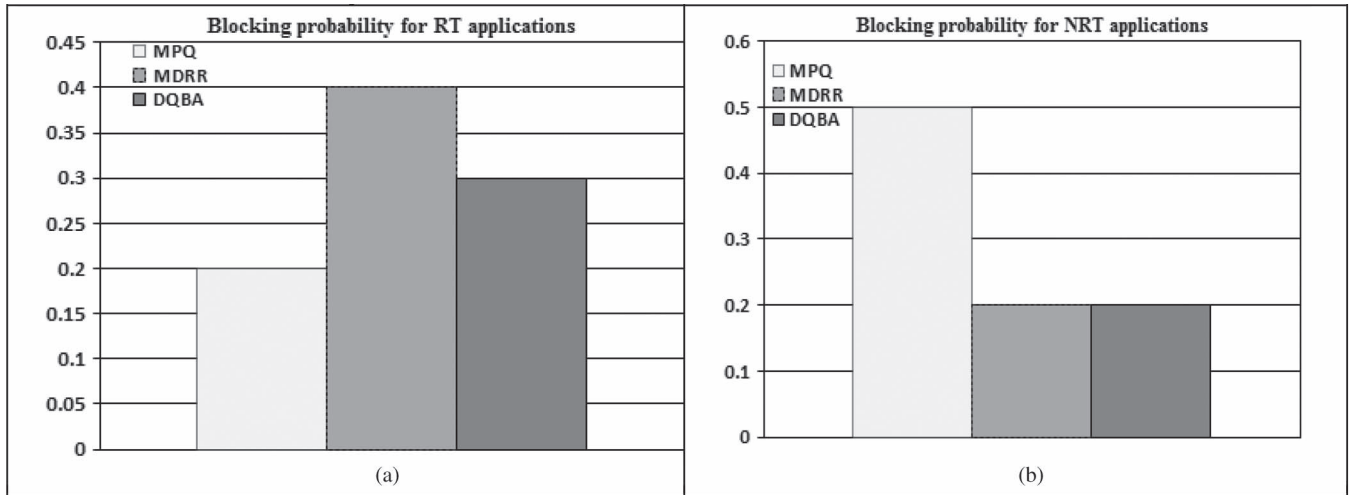


Fig. 15. (a) Blocking probability for RT applications using three models. (b) Blocking probability for NRT applications using three models.

and their QoS requirements are satisfied using our modified version of RR that takes into account fairness and utilization in the Tier-1 component of the proposed framework.

This figure shows the delay comparison for nrtPS flows in an intraclass level. The delay values increase to levels above 300 ms and then settle down at around 300, satisfying the typical benchmark values presented in Table I. This figure also shows that DQBA could deliver QoS guarantees at variable levels for several NRT flows, which are separated at the intraclass level.

### C. CAC Results

The CAC policy is designed to regulate the admission of the incoming calls, check the feasibility of accepting the calls based on available resources, and decide whether to accept or reject the calls. In this paper, the criteria for accepting the calls are based on three factors, namely, available bandwidth, fairness, and utilization. If a call passes all the tests, then it is accepted; otherwise, it is rejected.

Here, we ran the experiments based on the CAC blocking probability for all the incoming calls that are rejected by CAC using the three models and collected data on the number of calls received or blocked in each case. The results are presented in Fig. 15 for RT and NRT applications using all three models.

In the case of RT applications, MPQ blocked 20% of the calls based on fairness and utilization tests, whereas MDRR and DQBA blocked 40% and 30%, respectively. On the other hand, for NRT applications, where MPQ has blocked half of the calls, MDRR and DQBA both managed to allow 80% of the calls.

Fig. 15(a) shows that MPQ blocks the lowest portion of the RT service flows in comparison with other methods. This is due to the fact that MPQ provides higher levels for fairness and utilization of RT applications. MDRR, on the other hand, blocks twice that value. MDRR provides lower values for both fairness and utilization of RT applications. Therefore, it is expected to see higher rates of blocking probabilities. In terms of RT service calls, DQBA drops a moderate level at 30%. This is due to the nature of the proposed framework since it has the highest level of utilization.

In terms of NRT applications, MPQ shows the highest level of blocking probability. This reflects the typical nature of MPQ since it is not fair to NRT applications. Therefore, it provides lower BA and high values for utilization of this type of applications. As a consequence, the call blocking rate for NRT applications has increased to almost 50%. MDRR, on the other hand, while providing higher values for both fairness and utilization of NRT applications, has the lowest rate of call blocking. The proposed DQBA shows higher values close to 1.0 in both RT and NRT applications for Fr and Ut. One of the main goals of the proposed solution is to increase utilization, and this directly affects the reduction of the call blocking probability, which is apparent in Fig. 15. The proposed solution matches the values of call dropping probabilities for those of MDRR.

In the future, we will further investigate rate-related resource allocation and develop equations to relate the dCap values to the initial traffic arrival rate. We will also perform more experiments to compare “rate-related” and “bandwidth-request-related” BA schemes.

QoS differentiation is characterized by the ability of the RRM scheme in guaranteeing fairness to all service flows simultaneously, by improving utilization of the system. QoS differentiation for WiMAX has been recently studied in the literature, but there is neither a clear definition nor a set of metrics to quantify differentiation. The main motivation for the future study is to open a direction in the study of QoS support, which yields to further quantify and analyze QoS differentiation using the proposed metrics.

We will also perform further study on fairness and utilization in the context of QoS differentiation and carry out a comparative analysis for resource allocations using traffic behavior and network conditions, as opposed to resource allocation that is only based on resource request. We plan to further develop this scheme by incorporating new metrics into the core of both intra- and interclass QoS support. This will allow us to further investigate the performance of various parts of the framework using QoS differentiation ideas and perform a comprehensive study on “quantifying QoS differentiation” such as variations in fairness and utilization.

## V. CONCLUSION AND FUTURE WORK

In this paper, we have proposed a framework, which we call DQBA, for the RRM techniques in WiMAX, which supports all types of service flows and makes BA that is dynamic, fair, and efficiently utilized. We showed by simulation results that DQBA could deliver QoS support while being fair to all classes of service defined by the standard. We also introduced two new metrics to evaluate QoS differentiation in WiMAX, namely, fairness and utilization. We used the new metrics to perform PS, CAC, and BA using traffic behavior and network conditions.

DQBA shows superior performance in the following two cases: 1) in comparison with the MPQ model for NRT applications and 2) in comparison with MDRR for RT applications. However, it provides compatible performance in both cases with respect to RT applications using MPQ and with respect to NRT applications using MDRR. In terms of the overall performance of the system with respect to both RT and NRT applications, DQBA clearly outperforms both MPQ and MDRR models.

In addition, DQBA shows a granular level of traffic separation in both inter- and interclass levels and QoS support guarantees to both RT and NRT applications in both levels. In terms of call blocking probabilities for RT applications, the MPQ model rejects the least amount of calls, with DQBA and MDRR following. In terms of NRT applications, MPQ rejects 50% of the calls, whereas both DQBA and MDRR allow 80% of the calls.

## REFERENCES

- [1] A. Sayenko, O. Alanen, J. Karhula, and T. Hamalainen, "Ensuring the QoS requirements in IEEE 802.16 scheduling," in *Proc. 9th ACM Int. Symp. Model. Anal. Simul. Wireless Mobile Syst.*, Malaga, Spain, 2006, pp. 108–117.
- [2] IEEE Standard for Local and Metropolitan Area Networks, Air Interface for Fixed Broadband Wireless Access Systems, Amendment2: Physical and Medium Access Control Layers for Combined Fixed and Mobile Operation in Licensed Bands and Corrigendum 1, IEEE Std. 802.16e, Feb. 2006.
- [3] IEEE Standard for Local and Metropolitan Area Networks—Part 16: Air Interface for Fixed Broadband Wireless Access Systems, IEEE Std. 802.16d-2004, 2004.
- [4] F. Hou, P. Ho, and X. Shen, "Performance evaluation for unsolicited grant service flows in IEEE 802.16 networks," in *Proc. IWCNC*, Vancouver, BC, Canada, 2006, pp. 991–996.
- [5] J. G. Andrews, A. Ghosh, and R. Muhamed, *Fundamentals of WiMAX: Understanding Broadband Wireless Networking*. Englewood Cliffs, NJ: Prentice-Hall, 2007.
- [6] J. Chen, W. Jiao, and Q. Guo, "Providing integrated QoS control for IEEE 802.16 broadband wireless access systems," in *Proc. VTC*, Dallas, TX, 2005, vol. 2, pp. 1254–1258.
- [7] R. Jayapavathy and G. Sureshkumar, "Performance evaluation of scheduling schemes for fixed broadband wireless access systems," in *Proc. 13th IEEE ICN*, Boston, MA, 2005, pp. 2–6.
- [8] K. Wongthavarawat and A. Ganz, "Packet scheduling for QoS support in IEEE 802.16 broadband wireless access systems," *Int. J. Commun. Syst.*, vol. 16, no. 1, pp. 81–96, Feb. 2003.
- [9] J. Kim and K. Kim, "Apparatus and method for download packet scheduling in base station of a portable Internet system," U.S. Patent Application No: US 2007/0121636 A1, May 31, 2007.
- [10] M. Mehrjoo, M. Dianati, X. Shen, and K. Naik, "Opportunistic fair scheduling for the downlink of IEEE 802.16 wireless metropolitan area networks," in *Proc. 3rd Int. Conf. Quality Service Heterogeneous Wired/Wireless Netw.*, New York, 2006, vol. 191.
- [11] K. Vinay, N. Sreenivasulul, D. Jayaraml, and D. Das, "Performance evaluation of end-to-end delay by hybrid scheduling algorithm for QoS in IEEE 802.16 network," in *Proc. Int. Conf. Wireless Opt. Commun. Netw. (IFIP)*, Bangalore, India, 2006, pp. 5–10.
- [12] J. Liang, X. Sun, and N. Kang, "Performance evaluation of an integrated and efficient uplink scheduler for WiMAX network," in *Proc. IEEE ICCTA*, Oct. 16–18, 2009, pp. 326–330.
- [13] X. Bai, A. Shami, and Y. Ye, "Robust QoS control for single carrier PMP mode IEEE 802.16 system," *IEEE Trans. Mobile Comput.*, vol. 7, no. 4, pp. 416–429, Apr. 2008.
- [14] A. Belghith, "Pricing-based schedulers for WiMAX," in *Proc. IEEE Int. Conf. Wireless Mobile Comput., Netw. Commun., WIMOB*, Oct. 12–14, 2009, pp. 202–207.
- [15] S. Pizzi, A. Molinaro, and A. Iera, "On the performance of compensation based and greedy scheduling for the IEEE 802.16 standard," in *Proc. IEEE Int. Conf. Commun.*, Dresden, Germany, Jun. 2009, pp. 978–984.
- [16] S. Pizzi, A. Molinaro, and A. Iera, "AMC and channel-awareness for QoS-based scheduler design in WiMAX networks," in *Proc. Eur. Wireless Conf. (EW)*, Apr. 12–15, 2010, pp. 857–864.
- [17] M. A. S. Khan, A. Sattar, T. Mustafa, and S. Ahmad, "Performance evaluation and enhancement of uplink scheduling algorithms in point to multipoint WiMAX networks," *Eur. J. Sci. Res.*, vol. 42, no. 3, pp. 491–506, Jun. 2010.
- [18] C. F. Ball, F. Trembl, X. Gaube, and A. Klein, "Performance analysis of temporary removal scheduling applied to mobile WiMAX scenarios in tight frequency reuse," in *Proc. 16th Annu. IEEE Int. Symp. PIMRC*, Berlin, Germany, Sep. 2005, pp. 888–894.
- [19] M. Dianati, X. Shen, and S. Naik, "A new fairness index for radio resource allocation in wireless networks," in *Proc. IEEE WCNC*, 2005, pp. 712–717.
- [20] B. Kaarthick, N. Nagarajan, A. Raja Mohamed, and G. Saimethun, "CINR and  $n$ -factor dependent fair scheduling algorithm for Mobile WiMAX," *Int. J. Recent Trends Eng.*, vol. 2, no. 1, pp. 63–69, Nov. 2009.
- [21] L. Nuaymi and Z. Noun, "Simple capacity estimation in WiMAX/802.16 system," in *Proc. 17th Annu. IEEE Int. Symp. PIMRC*, Summer 2009, pp. 1–5.
- [22] D. Tarchi, R. Fantacci, and M. Bardazzi, "Quality of service management in IEEE 802.16 wireless metropolitan area networks," in *Proc. IEEE ICC*, Istanbul, Turkey, 2006, vol. 4, pp. 1789–1794.
- [23] A. Zahran, B. Liang, and A. Saleh, "Beyond 3G wireless network design for optimal resource utilization," in *Proc. 23rd Biennial Symp. Commun.*, 2009, pp. 256–260.
- [24] D. Chuck and J. M. Chang, "Bandwidth recycling in IEEE 802.16 networks," *IEEE Trans. Mobile Comput.*, vol. 9, no. 10, pp. 1451–1464, Oct. 2010.
- [25] Cisco IOS Quality of Service Solutions Configuration Guide, Rel. 12.1, 2011. [Online]. Available: [http://www.cisco.com/en/US/docs/ios/12\\_1/qos/configuration/guide/qcdconmg.html#wp1001203](http://www.cisco.com/en/US/docs/ios/12_1/qos/configuration/guide/qcdconmg.html#wp1001203)
- [26] W. Odom and M. J. Cavanaugh, *Deploying Quality of Service in the Enterprise Networks*. San Jose, CA: Cisco, 2006.
- [27] H. Rashwan, H. M. ElBadawy, and H. H. Ali, "Comparative assessments for different WiMAX scheduling algorithms," in *Proc. WCECS*, San Francisco, CA, Oct. 2009, vol. I.



**Amir Esmailpour** received the Ph.D. degree from the University of Guelph, Guelph, ON, Canada.

He was a Software/Network Engineer with the telecom sector for several years, prior to his return to research and academia. He is currently a Post-doctoral Fellow with the Department of Electrical and Computer Engineering, University of Toronto, Toronto, ON. He has authored several papers in reputable publications such as the IEEE journals and conferences.



**Nidal Nasser** (M'00) received the Ph.D. degree from the School of Computing, Queen's University, Kingston, ON, Canada in 2004.

He is currently an Associate Professor with the School of Computer Science, University of Guelph, Guelph, ON.

Dr. Nasser is a member of several IEEE technical committees. He is an Associate Editor of the *Journal of Computer Systems, Networks, and Communications*, Wiley's *International Journal of Wireless Communications and Mobile Computing*, and Wiley's *Security and Communication Networks Journal*.