

Dynamic Bandwidth Request-Allocation Algorithm for Real-Time Services in IEEE 802.16 Broadband Wireless Access Networks

Eun-Chan Park*, Hwangnam Kim[†], Jae-Young Kim*, and Han-Seok Kim*

* Telecommunication & Network Division, Samsung Electronics Co., LTD., Korea

E-mail: {eunchan.park|jay.m.kim|hs365.kim}@samsung.com

[†] School of Electrical Engineering, Korea University, Korea, E-mail: hnkim@korea.ac.kr

Abstract—The emerging broadband wireless access (BWA) technology based on IEEE 802.16 is one of the most promising solutions to provide ubiquitous wireless access to the broadband service at low cost. This paper proposes an efficient uplink bandwidth request-allocation algorithm for variable-rate real-time services in IEEE 802.16 BWA networks. In order to minimize bandwidth wastage without degrading quality of service (QoS), we introduce a notion of *target delay* and propose *dual feedback architecture*. The proposed algorithm calculates the amount of bandwidth request such that the delay is regulated around the desired level to minimize delay violation and delay jitter for real-time services. Also, it can maximize utilization of wireless channel by making use of dual feedback, where the bandwidth request is adjusted based on the information about the backlogged amount of traffic in the queue and the rate mismatch between packet arrival and service rates. Due to the dual feedback architecture, the proposed scheme responds quickly to the variation of traffic load and is robust to the change of network condition. We analyze the stability of the proposed algorithm from a control-theoretic viewpoint and derive a simple design guideline based on the analysis. By implementing the algorithm in *OPNET* simulator, we evaluate its performance in terms of queue regulation, optimal bandwidth allocation, delay controllability, and robustness to traffic characteristics.

Index Terms—IEEE 802.16, uplink scheduling, bandwidth request, quality of service, real-time service

I. INTRODUCTION

In recent years, broadband wireless access (BWA) networks have been rapidly evolved to satisfy increasing demands of users for ubiquitous and seamless access to the broadband service, such as video conferencing, real-time multimedia streaming, Internet Protocol TV, as well as traditional Internet services under mobile wireless environments. The emerging IEEE 802.16e BWA network [1], called *mobile WiMAX*, is one of the most promising solutions for the last mile broadband wireless access to support high data rate, high mobility, and wide coverage at low cost. In 2006, the *Wireless Broadband* (WiBro) service, the first commercial service of mobile WiMAX over the world, was launched in South Korea, and Sprint declared to provide mobile WiMAX service in united state from April, 2008. With the rapid growth of real-time and/or multimedia service, providing quality of service (QoS) in BWA networks is an imperative and challenging issue.

In order to support QoS for various types of traffic, IEEE 802.16 medium access control (MAC) protocol [1] defines several bandwidth request-allocation mechanisms and five types of scheduling classes; Unsolicited Grant Service (UGS), real-time Polling Service (rtPS), non-real-time Polling Service

(nrtPS), Best-Effort (BE), and extended real-time Polling Service (ertPS). Both UGS and rtPS are proposed to support real-time service generating packets periodically. While UGS is suitable for constant bit rate (CBR) traffic such as Voice over Internet Protocol (VoIP), rtPS is for variable bit rate (VBR) traffic such as MPEG video. The UGS scheduling mechanism can minimize delay in bandwidth request-allocation process, however; at the same time, it might waste bandwidth or suffer from insufficient bandwidth with VBR traffic. On the other hand, the rtPS mechanism effectively utilizes bandwidth at the cost of additional delay due to on-demand bandwidth request. To strike a balance between minimizing delay and maximizing utilization, ertPS is introduced in IEEE 802.16e [1], an amendment of IEEE 802.16-2004 [2]. Like UGS, the ertPS scheduling mechanism allocates bandwidth periodically without any request so as to minimize delay. Also, in a similar way to rtPS, it can adjust the size of bandwidth allocation to maximize utilization. However, any specific bandwidth request-allocation algorithm is not standardized so that proprietary implementations may be used by equipment vendors. Although there have been several proposals for QoS scheduling frameworks and algorithms in IEEE 802.16 BWA networks in the literature [3]–[8], they mainly focus on the QoS architecture and scheduling algorithm in a base station to satisfy diverse QoS requirements, rather than bandwidth request algorithm in a subscriber station.

In this paper, we propose a simple and efficient uplink bandwidth request algorithm for the ertPS scheduling mechanism, aiming to minimize bandwidth wastage without violating QoS, i.e., to devise a dynamic bandwidth request-allocation mechanism. The key idea for this algorithm is twofold; (i) In order to maintain satisfactory QoS, we introduce a notion of *target delay*, i.e., tolerable delay in the MAC layer, which can be interpreted into the target value of the transmission queue length. (ii) Moreover to maximize utilization, we propose to deploy a *dual feedback* architecture; one for difference between the backlogged amount of traffic in the transmission queue and its target value and the other for mismatch between packet¹ arrival and service rates.

Using the dual feedback, the proposed algorithm dynamically calculates the amount of bandwidth request so that the bandwidth wastage is minimized. At the same time, it can

¹In this paper, a packet denotes a MAC-layer payload without including MAC-layer overhead, i.e., MAC service data unit (MSDU), unless otherwise stated.

minimize delay violation and delay jitter, by controlling the MAC-layer service delay around the desired level. Moreover, it responds quickly to the variation of traffic load and is robust to the change of network condition, due to the dual feedback architecture. Based on a control-theoretic approach, we analyze the performance and stability of the proposed algorithm and derive a simple design guideline. Also, we implement the algorithm in IEEE 802.16 MAC layer using OPNET [9] simulator, and perform extensive simulations. The simulation results confirm that the proposed algorithm can minimize bandwidth wastage and regulate delay around the desired level with significantly reduced jitter. In this paper, we restrict our attention on the uplink bandwidth request mechanism for VBR traffic, since downlink scheduling does not involve any bandwidth request-allocation process and adjusting the size of bandwidth request is not necessary for CBR traffic.

The rest of this paper is organized as follows. In Section II, we briefly introduce QoS scheduling architecture and uplink bandwidth request-allocation mechanisms standardized in IEEE 802.16e. In Section III, we propose a dynamic bandwidth request mechanism, which makes use of target delay and dual feedback architecture. In Section IV, we analyze the proposed algorithm and provide a design guideline to make the system stable. Section V presents extensive simulation results to evaluate the performance of the proposed scheme. Finally, we conclude the paper in Section VI.

II. QOS ARCHITECTURE OF IEEE 802.16 NETWORKS

A. Scheduling framework

This paper considers point-to-multipoint (PMP) architecture of IEEE 802.16 BWA networks, where transmission only occurs between a base station (BS) and subscriber stations (SSs). The BS controls all the communications between BS and SSs. All the transmissions are associated with a unidirectional connection, which is associated with a service flow characterized by a set of QoS parameters, e.g., tolerable delay and minimum/maximum traffic rate. The connection can be either downlink (from BS to SS) or uplink (from SS to BS), each of which is denoted as DL and UL, respectively. When establishing a connection, a proper connection admission control is performed at the BS. Once the connection is admitted, the scheduler in the BS schedules the DL and UL connections independently. Also, they are served in the separate region of physical (PHY) layer frame, e.g., orthogonal frequency division multiple access with time division duplex (OFDMA/TDD) frame. The DL channel is a broadcast channel, while the UL channel is shared by several SSs in a manner that a SS requests its required bandwidth and the BS allocates it by scheduling all the requests from the SSs. The scheduler in the BS generates and broadcasts MAP message containing two dimensional (time and frequency) channel allocation information for DL and UL connections. The UL MAP message specifies the time when a SS can transmit and how long it can do, and which sub-channel it can occupy.

Depending on the scheduling class, there are several ways of bandwidth request;

- (i) without any request from the SS, the BS allocates bandwidth periodically,
- (ii) the SS receives a periodic bandwidth request opportunity from the BS, poll,
- (iii) the SS contends for the bandwidth request opportunity.

The details about uplink bandwidth request-allocation mechanisms defined in the IEEE 802.16e will be discussed in the next subsection.

B. Uplink bandwidth request-allocation mechanisms

The standard of IEEE 802.16-2004 [2] defines four uplink scheduling classes;

- **UGS**: This class has the highest service priority and is designed to support CBR traffic, e.g., VoIP traffic. On establishing UGS connection, the SS declares its bandwidth requirement and maximum tolerable delay. Then, BS allocates the requested amount of bandwidth periodically in an unsolicited way. Therefore, UGS can eliminate overhead and delay resulting from the bandwidth request-allocation process. It is suitable for the applications requiring constant bandwidth allocation with minimal delay and jitter.
- **rtPS**: This is for real-time VBR traffic generating variable-sized packets periodically, e.g., MPEG video. By issuing polls at every given interval, the BS gives request opportunities to SSs. Then, the SS requests bandwidth without contending with other SSs. While UGS is proactive to the bandwidth requirement, rtPS is reactive to the bandwidth demand. Therefore, rtPS involves an additional delay in the bandwidth request-allocation process.
- **nrtPS**: This scheduling class is designed to support non-real-time VBR traffic that requires minimum bandwidth guarantee but is insensitive to delay, e.g., FTP. The nrtPS scheduling class uses the same polling mechanism as rtPS, however; it is allowed to contend for non-periodical bandwidth request opportunity.
- **BE**: This is for the best effort traffic that does not have any specific QoS requirements, e.g., e-mail or web. The BS does not give any dedicated request opportunity to the SSs, and the SS sends bandwidth request message in a contention-based way.

In addition to these four types of service classes, IEEE 802.16e [1] introduces another service class, ertPS.

- **ertPS**: This is basically identical to UGS, except that ertPS can change the amount of bandwidth allocation dynamically depending on the traffic characteristics. On detecting that the allocated bandwidth is insufficient to serve packets in time, the SS requests an additional bandwidth by piggybacking its amount on the packet header. Otherwise if the allocated bandwidth is excessive, the SS can request decreasing the amount of bandwidth allocation. Therefore, the ertPS is suitable for real-time VBR traffic and VoIP traffic with silence suppression.

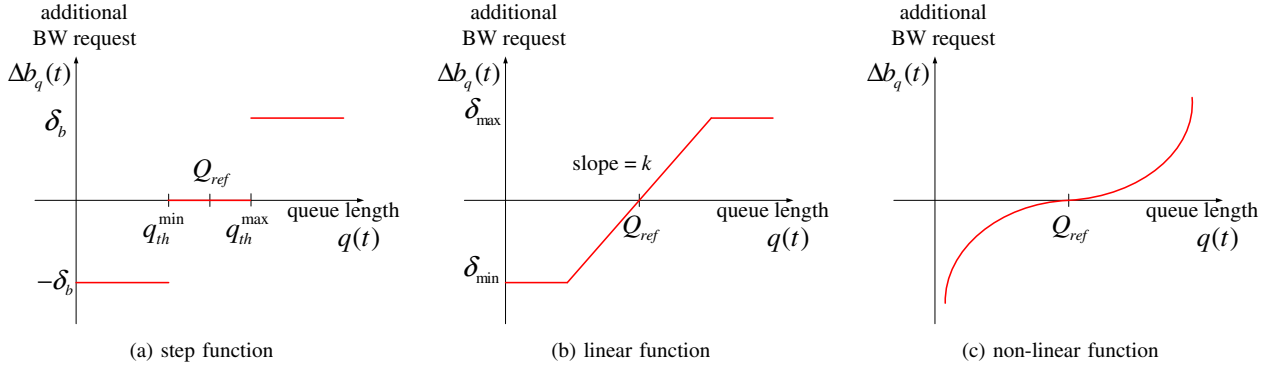


Fig. 1. Several functions to calculate additional bandwidth request depending on the queue length.

In summary, the bandwidth request for UGS is done in an unsolicited manner and that for rtPS and nrtPS is achieved in a polling-based way, and the BE service contends for bandwidth request opportunity. On the other hand, the bandwidth allocation for UGS is performed based on reservation and that for rtPS, nrtPS, and BE is done on a demand-basis. The ertPS class employs hybrid approach in bandwidth request and allocation.

III. DYNAMIC BANDWIDTH REQUEST ALGORITHM

A. Design rationale

The design objectives of the dynamic bandwidth request-allocation algorithm for ertPS are as follows;

- It should estimate the required bandwidth timely and accurately.
- It should neither waste bandwidth nor suffer from lack of bandwidth.
- It should minimize violation of delay requirement.

To achieve these purposes, we introduce **target delay** and **dual feedback**.

1) *Introduction of target delay*: The target delay plays a key role in determining the amount of bandwidth adjustment. Most real-time services have a tolerable end-to-end delay up to which QoS is not much deteriorated, e.g., 100 ~ 200 ms for VoIP service and a few hundreds of milliseconds for streaming service. Taking this tolerable delay into account, we can set a target MAC-to-MAC delay, T_{ref} (sec), and translate it into the target length of transmission queue, Q_{ref} (byte), under the following reasonable assumptions. (A1) A connection admission control is applied to the real-time services so that their total amount of bandwidth request does not exceed the available channel capacity on average. (A2) The scheduler deploys strict priority scheduling algorithm so that it allocates the requested bandwidth for the real-time service with high priority. (A3) An ertPS connection is established such that its bandwidth allocation interval, denoted as T_a (sec), is equal to the packetization interval of the real-time service, which is constant. (A4) The backhaul capacity of link connecting wireless access network to wired networks is large enough not to incur any queuing delay in the BS for UL connection. With these assumptions, Q_{ref} can be represented in terms of

T_{ref} and T_a as;

$$Q_{ref} = \frac{T_{ref} - T_o}{T_a} \bar{l}, \quad (1)$$

where \bar{l} (byte) denotes the average packet size and T_o (sec) represents any additional delay except queuing delay, e.g., processing delay at the MAC layer and transmission delay over wireless channel. It is noteworthy that T_{ref} does not include codec delay, packetization delay, and playout buffer delay at the application layer. Let us denote the size of transmission queue as $q(t)$ (byte) and the required additional bandwidth due to $q(t)$ as $\Delta b_q(t)$ (byte/sec). As $q(t)$ increases over Q_{ref} , $\Delta b_q(t)$ needs to be increased accordingly to satisfy delay requirement.

Now, we consider several approaches to calculate $\Delta b_q(t)$. Let us define δ_b as the amount of bandwidth required to transmit one packet whose size is \bar{l} . As shown in Fig. 1(a), the primitive way to control bandwidth request is to ask for δ_b additionally if $q(t)$ exceeds a maximum threshold value, q_{th}^{max} ($> Q_{ref}$), and to reduce the allocated bandwidth by amount of δ_b if $q(t)$ falls below a minimum threshold value, q_{th}^{min} ($< Q_{ref}$). Alternatively, we can calculate $\Delta b_q(t)$ in proportion to the difference between $q(t)$ and Q_{ref} , as depicted in Fig. 1(b). In this approach, we can set the upper and lower limit on $\Delta b_q(t)$, denoted as δ_{max} (> 0) and δ_{min} (< 0), respectively, to avoid an abrupt change of $\Delta b_q(t)$. Also, we can consider a non-linear function to calculate $\Delta b_q(t)$, as shown in Fig. 1(c). In this approach, $\Delta b_q(t)$ changes sharply as the discrepancy between $q(t)$ and Q_{ref} increases.

2) *Dual feedback approach*: The bandwidth request control based on the queue length, $\Delta b_q(t)$, reacts slowly to the variation of packet arrival rate because $\Delta b_q(t)$ changes after detecting the deviation of queue length from the desired level. To make the response fast, we introduce the dual feedback consisting of two feedback loops for queue length and rate. Let us define packet arrival rate and service rate as $a(t)$ and $s(t)$ (byte/sec), respectively, and an additional bandwidth request due to rate mismatch as $\Delta b_r(t)$ (byte/sec). As the information of queue length mismatch, $e_q(t) = q(t) - Q_{ref}$, is used in calculating $\Delta b_q(t)$, the information of rate mismatch, $e_r(t) = a(t) - s(t)$, is utilized in calculating $\Delta b_r(t)$. As the

packet arrival rate exceeds the service rate, i.e., $a(t) > s(t)$, packets start to be accumulated. In this case, $\Delta b_r(t)$ need to be positive to serve these packets timely without violating delay requirements. On the other hand, if $a(t) < s(t)$, the queue length tends to decrease. In this case, less bandwidth is required and $\Delta b_r(t)$ become negative in order not to waste bandwidth. The rate feedback provides predictive information about queue length. Consequently, the bandwidth request control based on the rate feedback shows anticipatory response to the queue length change, giving fast response to the variation of packet arrival rate. The total additional bandwidth request under the dual feedback architecture, $\Delta B(t)$, consists of queue-based component $\Delta b_q(t)$ and rate-based component $\Delta b_r(t)$ and can be represented in a generalized form as;

$$\begin{aligned}\Delta B(t) &= \Delta b_q(t) + \Delta b_r(t) \\ &= f(e_q(t)) + g(e_r(t)),\end{aligned}\quad (2)$$

where $f(\cdot)$ and $g(\cdot)$ indicate an appropriate non-negative function, as illustrated in Fig. 1.

B. Algorithm and implementation issues

In this subsection, we provide a detailed algorithm for the bandwidth request algorithm and discuss several issues regarding its implementation and overhead. From (2), we consider linear functions for $f(\cdot)$ and $g(\cdot)$ for the simplicity as

$$\Delta B(t) = K_q e_q(t) + K_r e_r(t), \quad (3)$$

where K_q and K_r denote constant control parameters that have non-negative values. In (3), the rate mismatch $e_r(t)$ can be represented in terms of queue length mismatch, i.e.,

$$\begin{aligned}e_r(t) &= a(t) - s(t) \\ &= \frac{d}{dt}q(t) = \frac{d}{dt}e_q(t), \quad \text{for } 0 < q(t) < Q_{max}.\end{aligned}\quad (4)$$

Here, Q_{max} is the maximum queue size. To implement this algorithm, we need to transform the continuous-time function $\Delta B(t)$ to the discrete-time function by sampling every bandwidth allocation interval, i.e., $\Delta B[n] = \Delta B(nT_a)$, where n is a non-negative integer. We approximate the derivative term in (4) using a first-order Euler approximation, i.e.,

$$\frac{d}{dt}e_q(t) \approx \frac{e_q[n] - e_q[n-1]}{T_a}. \quad (5)$$

From (3)–(5), we can calculate the bandwidth request $\Delta B[n]$ only using the current and previous values of queue length error. Remind that $\Delta B[n]$ is the increment or decrement of bandwidth request during the n th allocation interval and the corresponding total bandwidth request during this interval, $B[n]$, becomes

$$B[n] = \max(B[n-1] + \Delta B[n], B_{min}), \quad (6)$$

where B_{min} denotes the minimum amount of bandwidth allocation required for issuing bandwidth request.

After calculating the additional bandwidth $\Delta B[n]$, a SS informs the BS of $\Delta B[n]$, by conveying it on the *extended*

piggyback request (EPBR) field of *grant management sub-header* [1]. The size of EPBR field is 11 bits and it has two operation modes, incremental mode and aggregate mode. If the first bit of EPBR is set to zero, the remaining 10 bits represent increment of bandwidth request, otherwise they represent aggregate bandwidth request. Therefore, if $\Delta B[n] > 0$, it can be carried with both incremental and aggregate modes. Otherwise if $\Delta B[n] < 0$, the SS calculates $B[n]$ as given in (6) and carries this value with the aggregate mode.

In order to apply the proposed algorithm to the case of on-off traffic, e.g., VoIP traffic with silence suppression, we need to elaborate this algorithm. If the queue length maintains zero longer than a given threshold time, we determine that the connection becomes inactive, and stop the adaptation process for the bandwidth request. During this period, the total bandwidth request becomes its minimum value. If the connection becomes active again, which can be detected once the queue length becomes larger than zero, we restart the adaptation process with the initial value of $B[n]$, which is declared in the QoS parameters of the connection.

The communication overhead associated with piggybacking the bandwidth request is two bytes, the size of grant management subheader (i.e., $B_{min} = 2$ in (6)). Note that this minimal overhead is inevitable for the real-time scheduling services of IEEE 802.16 (e.g., UGS and rtPS, as well as ertPS). On the other hand, the computation overhead of the proposed algorithm is not significant. The number of operations is quite small and a SS only has to keep track of the current and previous values of its own queue length without estimating packet arrival and service rate explicitly (refer (5)). Moreover, the calculation for the bandwidth request is performed by each SS in a distributed manner. Thus, the proposed algorithm does not degrade scalability of the BS. Although the proposed dual feedback architecture is developed for the ertPS scheduling class, it can be applied to the rtPS class without any significant changes and it can be further extended to any request-based scheduling framework.

IV. ANALYSIS OF THE PROPOSED ALGORITHM

This section analyzes the behavior of the proposed algorithm. A system model is derived and its stability is analyzed from a control-theoretic viewpoint. Based on this analysis, a simple design guideline for control parameters is provided.

A. System modeling

We can model the bandwidth request-allocation mechanism described in Section III with three dynamic equations for (i) queue length error, (ii) additional bandwidth request, and (iii) total bandwidth request. For the sake of tractability, we consider continuous-time model, instead of discrete-time model. From (3), (4), and (6), the overall system can be regarded as a third-order linear feedback system with time delay;

$$\dot{e}_q(t) = a(t) - s(t), \quad (7)$$

$$\Delta \dot{B}(t) = K_q \dot{e}_q(t) + K_r \ddot{e}_q(t), \quad (8)$$

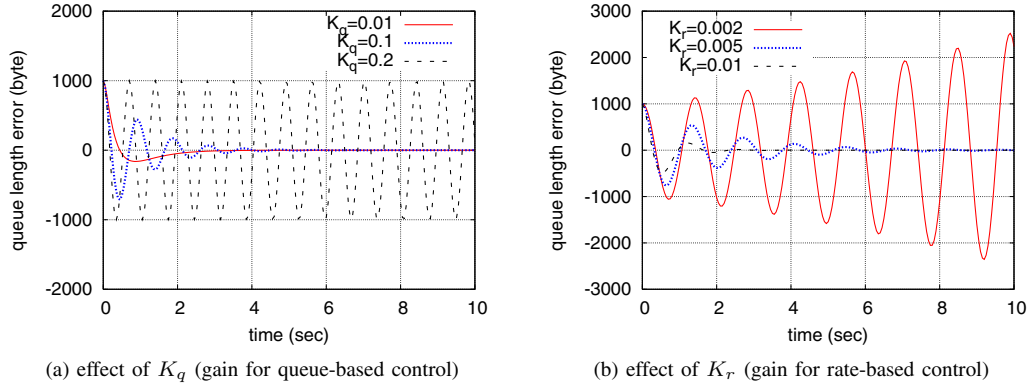


Fig. 2. Impulse response of the proposed system with various values of control parameters.

$$\dot{B}(t) = \frac{\Delta B(t)}{T_a}, \quad (9)$$

$$s(t) = \alpha \frac{B(t - T_a)}{T_a}, \quad (10)$$

The dynamic equation (9) is derived from (6) with the assumption that $B[n-1] + \Delta B[n] > B_{min}$. Here, (10) represents the model for the bandwidth allocation algorithm in the BS and $\alpha (\leq 1)$ denotes the ratio of bandwidth allocated by the BS to bandwidth requested by the SS. Under the assumptions of admission control and priority scheduling, i.e., (A1) and (A2), $\alpha \approx 1$ on average.

We investigate the behavior and stability of the proposed bandwidth request-allocation mechanism using the transfer function. Taking Laplace Transform to the system model described in (7) – (10) and approximating the time delay term as a first-order lag, i.e., $e^{-sT_a} \approx 1/(1 + T_a s)$, we have the following transfer function

$$G(s) = \frac{E_q(s)}{A(s)} = \frac{s^2 + \frac{1}{T_a}s}{s^3 + \frac{1}{T_a}s^2 + \frac{\alpha}{T_a^3}(K_r s + K_q)}. \quad (11)$$

The transfer function (11) is characterized by two control parameters (K_q and K_r) and the bandwidth allocation interval T_a .

B. Effect of control parameters

In this subsection, we investigate the behavior of the proposed algorithm using numerical analysis. For the simplicity of analysis, we assume that the maximum queue size is infinite and set $\alpha = 1$. Figures 2(a) and 2(b) show impulse responses of the proposed system with several values of K_q and K_r , respectively. The impulse response shows how $e_q(t)$ decays and is stabilized for the impulse input of $a(t)$. The default values for K_q and K_r are set to 0.05 and 0.01, respectively.² First, we investigate the effect of K_q from Fig. 2(a). If $K_q = 0.01$, $e_q(t)$ converges to zero with negligible fluctuation. However, the fluctuation of $e_q(t)$ increases as the value of K_q increases. For some critical value of K_q , e.g., 0.2 in

this numerical analysis, $e_q(t)$ oscillates continuously without converging. Moreover, we observed that it diverges as long as K_q exceeds this critical value. Next, we observe the effect of K_r on system stability from Fig. 2(b). In the case of $K_r = 0.002$, $e_q(t)$ oscillates continuously and diverges, which implies that the amount of bandwidth request can increase infinitely. However, if the value of K_r exceeds this critical value, the system becomes stable and the fluctuation of $e_q(t)$ decreases as K_r increases. By comparing Fig. 2(a) and Fig. 2(b), we can observe the followings:

- The system becomes unstable as K_q increases or K_r decreases.
- If the rate-based control for bandwidth request is disabled, i.e., $K_D = 0$, $e_q(t)$ diverges.

These analysis results support the importance of the dual feedback architecture conveying information of queue length error and rate mismatch.

C. Stability analysis

As confirmed in Fig. 2, the queue length converges to the desired target value depending on the control parameters. It is imperative to derive condition of control parameters that assures system stability. For this purpose, we provide the following stability criterion, which can be used as a design guideline for control parameters.

PROPOSITION. *The proposed bandwidth request-allocation system employing the dual feedback architecture is stable, i.e., the queue length converges to the desired target value in the steady state, if and only if the control parameters satisfy the following criterion.*

$$K_r - T_a K_q > 0 \quad (12)$$

Proof: From (11), the characteristic equation is given as

$$s^3 + \frac{1}{T_a}s^2 + \frac{\alpha}{T_a^3}(K_r s + K_q) = 0. \quad (13)$$

Let us define the coefficient of n th-order term in (13) as a_n , e.g., $a_2 = 1/T_a$. Since the given system is a linear time-invariant system, the stability condition can be derived by applying the Routh-Hurwitz criterion [10] to (13). The

²The values of control parameters are set in a trial-and-error manner so as to illustrate their effects. Later, we will provide a design guideline based on stability analysis.

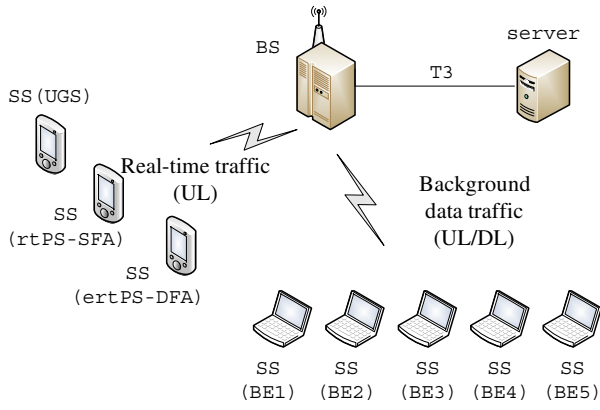


Fig. 3. Simple network configuration used in the simulations.

resulting condition becomes $a_1a_2 - a_3a_0 > 0$, which is identical to (12). ■

Corollary. *The bandwidth request-allocation system without rate-based control is unstable, regardless of the value of queue-based control gain. Therefore, the rate feedback is crucial for stability.*

Proof: It is obvious that the stability condition (12) cannot be satisfied with $K_r = 0$. ■

V. PERFORMANCE EVALUATION

In this section, we conduct extensive simulations using OPNET simulator with WiMAX module [9] to evaluate the performance of the proposed algorithm in several aspects and to compare it with other algorithms regarding bandwidth allocation efficiency and QoS assurance.

A. Simulation setup

We consider an IEEE 802.16 BWA network under PMP architecture consisting of multiple SSs and a BS, as shown in Fig. 3. The BS is connected to a server via a T3 link. The OFDMA parameters and their values are listed in Table I. The empirical COST-231 HATA model [11] is used for path-loss model, and ITU channel model [12] is used for multipath fading effect. Also, the shadowing is modeled as a lognormal random variable with zero mean and standard deviation of 8.9 dB. In the OPNET simulator, we emulated hybrid automatic repeat request mechanism so that the packet error rate is attained around 1%. Furthermore, we implemented adaptive modulation and coding scheme in the OPNET simulator. The supported modulation and coding rate schemes (MCSs) are as follows; QPSK (1/12, 1/8, 1/4, 1/2, 3/4), 16QAM (1/2, 3/4), 64QAM (2/3, 3/4, 5/6) for downlink and QPSK (1/12, 1/8, 1/4, 1/2, 3/4), 16QAM (1/2, 3/4) for uplink. The MCS changes dynamically depending on the signal to interference noise ratio, which is obtained via separate link level simulation. The maximum transmission queue size of SS is set to 64 Kbytes. Unless otherwise stated, the simulation configuration follows the methodology recommended by WiMAX forum [13].

The performance will be evaluated in terms of efficiency, delay, and jitter:

TABLE I
IEEE 802.16E OFDMA PHY PARAMETERS AND THEIR VALUES USED IN THE SIMULATIONS.

parameter	value
base frequency	2.5 GHz
channel bandwidth	10 MHz
TDD frame duration	5 ms
cyclic prefix duration	11.42 μ sec
basic symbol duration	91.43 μ sec
Fast Fourier Transform size	1024
number of symbols	29/18 (DL/UL)

- **MAC-layer efficiency:** defined as B_{tx}/B_{alloc} , where B_{tx} and B_{alloc} are total bytes of MSDU sent and total bytes of bandwidth allocation³ received during the whole simulation time, respectively.
- **MAC-to-MAC delay:** defined as the time difference between t_{tx} and t_{rx} , where t_{tx} is the instant at which a packet is delivered from the IP layer to the MAC layer at the sender (SS) and t_{rx} is the instant at which the received packet is delivered from the MAC layer to the IP layer at the receiver (BS).
- **jitter:** defined as standard deviation (STD) of MAC-to-MAC delay⁴.

We compare these performance indices for the following three bandwidth request-allocation mechanisms;

- **UGS:** This is a baseline algorithm that does not employ any adaptive bandwidth request mechanism. The fixed amount of bandwidth is allocated periodically to the SS.
- **rtPS-SFA:** At every polling interval, the SS requests bandwidth at the amount of backlogged traffic in the transmission queue. This is a simple bandwidth request mechanism with single feedback algorithm (SFA) based on queue length. Also, this mechanism does not take *target delay* into account.
- **ertPS-DFA:** This deploys the proposed bandwidth request mechanism (dual feedback algorithm (DFA) with *target delay*) with ertPS. The control parameters are set as $K_q = 0.05$ and $K_r = 0.02$.

The rtPS-SFA algorithm can be regarded as a special case of the ertPS-DFA algorithm, i.e., rtPS-SFA is equivalent to ertPS-DFA if $K_q = 1$, $K_r = 0$, and $Q_{ref} = 0$.

A VBR video traffic is used in the simulations. Its packet size is randomly distributed while the packetization interval is fixed. At the same time, we consider the background traffic with greedy FTP traffic. The real-time traffic is transferred using RTP/UDP/IP protocol suite, while background traffic uses TCP/IP protocol suite. The three SSs in Fig. 3 send uplink video traffic, each of which has UGS, rtPS, and ertPS connection, and the remaining five SSs send/receive background traffic with BE connections.

³MAC-layer overhead is not included in B_{alloc} .

⁴Alternatively, the jitter may be defined as the time difference between two consecutive packets arrived at the receiver. In this study, however, it is defined as the standard deviation of delay to evaluate the queue regulation performance.

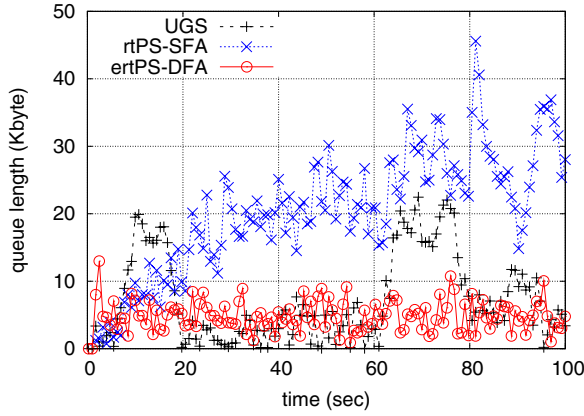


Fig. 4. Queue length for the UGS, rtPS-SFA, and ertPS-DFA algorithms.

B. Queue regulation and rate adaptation

In the first simulation, we evaluate the performance focusing on queue regulation and rate adaptation. The application-layer payload size is set to be exponentially distributed with mean of 500 bytes and packetization interval is set to 20 ms. Both bandwidth allocation interval of UGS and ertPS-DFA and polling interval of rtPS-SFA are set to be equal to the packetization interval.

Fig. 4 shows the queue lengths for the UGS, rtPS-SFA, and ertPS-DFA algorithms. In the case of UGS, $q(t)$ is minimized during some period, e.g., $t = 20\text{s} \sim 60\text{s}$, but it increases notably during some period, e.g., $t = 60\text{s} \sim 80\text{s}$. The reason is that UGS allocates fixed amount of bandwidth regardless of packet arrival rate. The queue length of rtPS-SFA fluctuates severely, causing large delay jitter. Also, $q(t)$ of rtPS-SFA is quite larger than those of UGS and ertPS-DFA, which results from the additional polling delay and bandwidth request only based on the current queue length. However, the ertPS-DFA algorithm regulates its queue length around the target value, which is set to 5 KB in this simulation, during the whole simulation time.

The reason of queue length variation of UGS and rtPS-SFA in Fig. 4 can be explained from the rate mismatch between the packet arrival and service rates shown in Fig. 5. In the case of UGS (Fig. 5(a)), its service rate cannot exceed the predefined value. In this simulation, the reserved bandwidth for UGS and ertPS is set to 220 Kb/s, by taking the packet overhead into consideration.⁵ Therefore, during some period (e.g., $t = 10\text{s} \sim 20\text{s}$ and $t = 60\text{s} \sim 70\text{s}$) where $a(t)$ exceeds $s(t)$ (see Fig. 5(a)), the queue length increases abruptly (see Fig. 4). Otherwise if $a(t) < s(t)$, the allocated bandwidth is wasted. On the other hand, rtPS-SFA requests bandwidth on demand, thus, it does not waste bandwidth. However, the bandwidth is not allocated timely, i.e., it is allocated only after request. Thus, as shown in Fig. 5(b), there occurs a small time shift between $a(t)$ and $s(t)$, while there is not a significant discrepancy between their magnitudes. For example, during $t = 80\text{s} \sim$

100s, we can observe a small but not negligible time delay between $a(t)$ and $s(t)$, causing a sudden large increase in $q(t)$ (see Fig. 4). In the case of ertPS-DFA, $s(t)$ agrees quite well with $a(t)$ as confirmed in Fig. 5(c). The ertPS-DFA algorithm makes use of rate feedback, as well as queue length feedback. The outstanding queue regulation performance of ertPS-DFA results from its rate adaptation. The queue regulation and rate adaptation of ertPS-DFA are key features to provide optimal bandwidth allocation and delay control.

C. Optimal bandwidth allocation

This simulation focuses on the performance from the perspective of optimal bandwidth allocation. It is clear that there is a trade-off between efficiency and performance in allocating bandwidth. As the allocated bandwidth increases, the performance (e.g., delay and jitter) is improved at the cost of decreased efficiency. In order to evaluate this trade-off for the various conditions of bandwidth allocation, we introduce *provisioning level*, ρ , defined as the ratio of the reserved bandwidth to the average packet arrival rate, and observe performance indices with respect to ρ ranging from 0.8 to 1.2.⁶ Unless otherwise stated, the simulation configuration is same as that in the previous simulation.

As shown in Fig. 6(a), the efficiency of UGS is nearly equal to 1 if $\rho < 1$, but it decreases almost linearly as ρ increases over 1. The bandwidth wastage of UGS exceeds 15% in the case of $\rho = 1.1$. However, both rtPS-SFA and ertPS-DFA maintain high level of efficiency close to 1 regardless of the value of ρ , as shown in Fig. 6(a). The efficiency of ertPS-DFA is slightly higher than that of rtPS-SFA. Next, we investigate the effect of ρ on delay from Fig. 6(b). Among three algorithms, the delay of UGS is most sensitive to the provisioning level. Once $\rho > 1$, i.e., excessive amount of bandwidth is allocated, the delay of UGS is minimized. Otherwise if $\rho < 1$, its delay is larger than that of ertPS-DFA up to about ten times.⁷ By comparing Fig. 6(a) and Fig. 6(b), we can observe that UGS has a clear trade-off between efficiency and delay with respect to the provisioning level. However, the delay of rtPS-SFA is not much affected by ρ and it decreases slightly as ρ increases. It is because rtPS-SFA requests bandwidth on demand every polling interval and its delay decreases as the polling interval decreases. The ertPS-DFA algorithm outperforms other algorithms significantly in terms of delay if $\rho < 1$. In this simulation, the target delay is set to 200 ms and the actual delay of ertPS-DFA lies between 214 ms and 222 ms for the entire range of ρ . Moreover, unlike other algorithms, the delay of ertPS-DFA is almost immune to the value of ρ .

⁶Accordingly, the reserved bandwidth of UGS and ertPS-DFA is set to range from $0.8 \times 216 \text{ Kb/s} \approx 175 \text{ Kb/s}$ to $1.2 \times 216 \text{ Kb/s} \approx 260 \text{ Kb/s}$. Also, ρ of rtPS-SFA is set by means of polling interval, ranging from $20 \text{ ms} / 0.8 = 25 \text{ ms}$ to $20 \text{ ms} / 1.2 = 16.7 \text{ ms}$.

⁷In order for the delay not to increase excessively, we can limit the size of transmission queue to a small value and apply *head drop* discipline, i.e., drop a packet in the head of transmission queue. This discipline may increase packet loss rate. In this simulation, the maximum queue size is set to a large value, 64 KB, to illustrate the effect of provisioning level on the delay without taking packet loss into account.

⁵The header sizes of RTP, UDP, and IP are 12, 8, and 20 bytes, respectively. Thus, the average arrival rate of MSDU with average application-layer payload size of 500 bytes, becomes $540 \times 8 \text{ bits} / 20 \text{ ms} = 216 \text{ Kb/s}$.

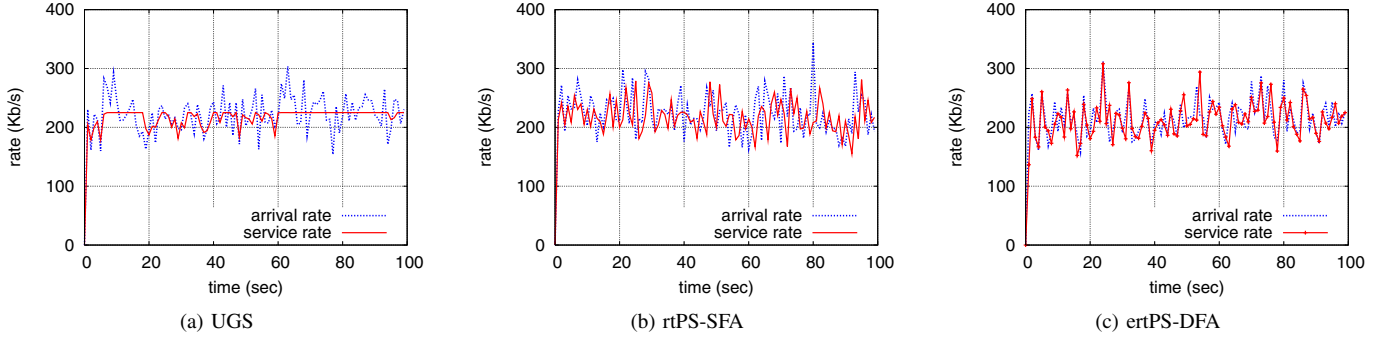


Fig. 5. Comparison of arrival rate and service rate for the UGS, rtPS-SFA, and ertPS-DFA algorithms.

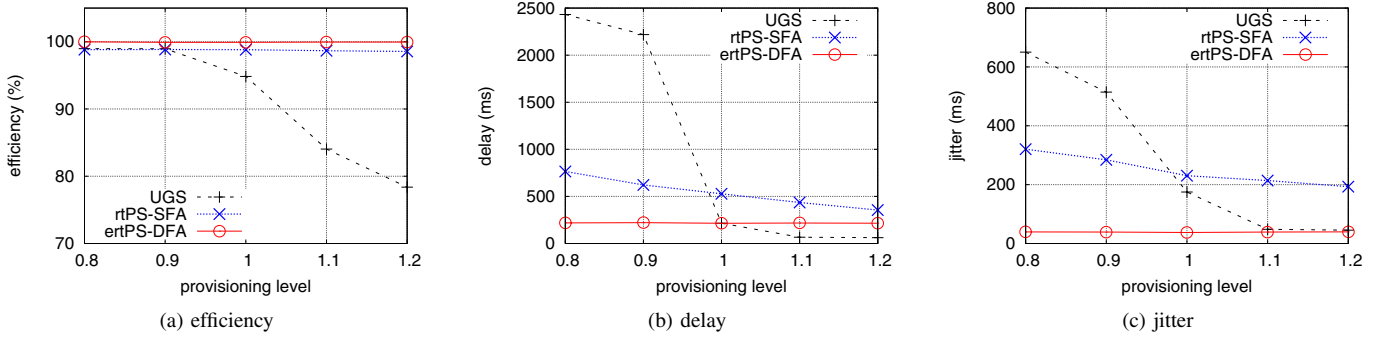


Fig. 6. Comparison of several performance indices with respect to the provisioning level.

The reason is that ertPS-DFA adaptively controls bandwidth request regardless of the predefined bandwidth reservation, so as to satisfy the desired target delay. We also compare the jitter of three algorithms in Fig. 6(c). As ρ increases, jitter of UGS decreases rapidly, but it does not decrease any more once $\rho > 1.1$. In the case of rtPS-SFA, jitter decreases slowly as ρ increases due to decrease of polling interval. The ertPS-DFA algorithm has almost constant delay jitter and its value is smaller than those of UGS and rtPS-SFA up to about 16 and 8 times, respectively, because ertPS-DFA tightly regulates queue length around the target value (see Fig. 4).

Fig. 6 confirms that ertPS-DFA allocates bandwidth in an optimal manner in that it satisfies delay requirement while minimizing bandwidth wastage and delay jitter, regardless of bandwidth provisioning level.

D. Delay control according to target value

We conduct this simulation to evaluate how the proposed ertPS-DFA algorithm can control delay according to the given target value. For this purpose, we repeat simulations with different values of target delay, 50, 100, 150, 200, and 300 ms.

Table II summarizes statistics of MAC-to-MAC delay, e.g., average, standard deviation, and 95% confidential interval. The difference between average delay and target delay is not larger than 27 ms, and it decreases as the target delay increases. The standard deviation is little affected by the target value, it ranges from 37.6 ms to 39.6 ms. Also, the range of 95% confidential interval is quite small, which confirms that the

TABLE II
DELAY STATISTICS OF THE ERTPS-DFA ALGORITHM WITH RESPECT TO THE SEVERAL TARGET VALUES.

target delay (ms)	average (ms)	STD (ms)	95% conf. interval (ms)
50	76.7	38.6	[75.6, 77.2]
100	119.4	37.6	[118.4, 120.5]
150	167.0	38.9	[166.0, 168.1]
200	216.1	39.6	[215.0, 217.2]
300	311.0	39.3	[309.9, 312.1]

delay is densely distributed around its average value. These results show the unique competence of the proposed ertPS-DFA algorithm in terms of fine controllability of delay, which cannot be achieved with UGS and rtPS-SFA algorithms.

E. Robustness to traffic characteristics

In this simulation, we investigate the effect of packet size and its distribution on the performance of the proposed ertPS-DFA algorithm. We change the average packet size \bar{l} from 100 bytes to 1000 bytes and consider the following three random distributions for the packet size;

- Exponential distribution: average = \bar{l} ,
- Uniform distribution: minimum = 1, maximum = $2\bar{l} - 1$,
- Pareto distribution: shape = 2, location = $\bar{l}/2$.

Among these distributions, the degree of randomness is weakest for the uniform distribution, and is strongest for the Pareto distribution. The packetization interval is set to 20 ms, thus, the application-layer bit rate changes from 40 kb/s to 400 kb/s,

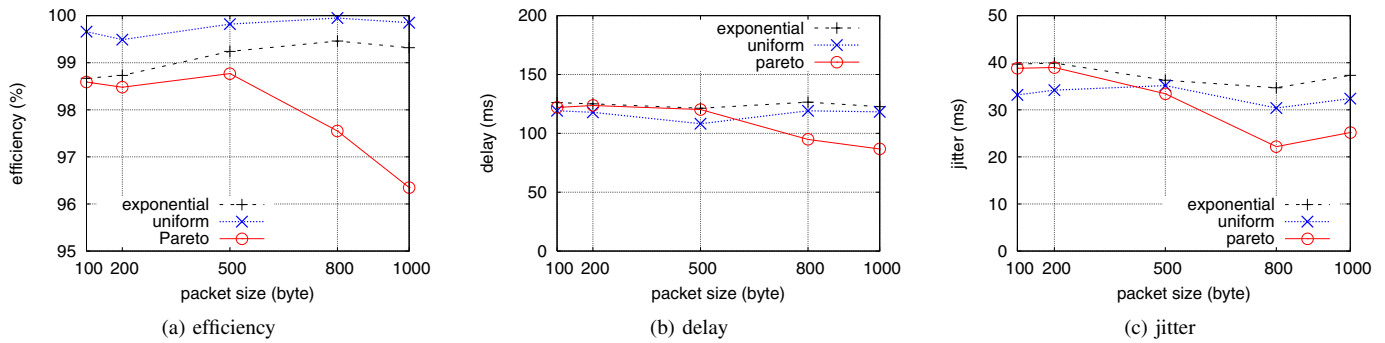


Fig. 7. Robust performance of the ertPS-DFA algorithm to the packet size and its distribution.

which is a reasonable configuration for VBR real-time traffic in IEEE 802.16 BWA networks.

As shown in Fig. 7(a), the efficiency of ertPS-DFA is best for the uniform distribution, it is higher than 99% for the entire range of \bar{l} . With the uniform and exponential distributions, \bar{l} has insignificant impact on the efficiency. In the case of Pareto distribution, however, the efficiency decreases notably as \bar{l} exceeds 500 bytes. Even for the worst case, the efficiency is larger than 96%. Next, we observe the delay and jitter from Fig. 7(b) and 7(c), respectively. The delay and jitter are not much affected by the packet size and its distribution. The target delay is set to 100 ms in this simulation, and the average delay lies between 87 ms and 126 ms. On the other hand, the jitter lies between 25 ms and 39 ms. The simulation results in Fig. 7 confirm the robust performance of the ertPS-DFA algorithm to various traffic patterns.

VI. CONCLUSION

We have proposed the dynamic bandwidth request mechanism for VBR real-time traffic in IEEE 802.16 broadband wireless access networks. By introducing the notion of target delay, a tolerable delay for real-time service, we can dynamically calculate the amount of bandwidth request that maximizes efficiency of wireless channel without violating delay requirement. In order to make the response to the change of traffic load fast, we have introduced the dual feedback architecture, where the difference between the actual queue length and the desired target length and the rate mismatch between packet arrival rate and service rate are utilized as feedback information. Due to the target delay and dual feedback architecture, the proposed algorithm tightly regulates the queue length around the desired level, thus it can control delay to the target level while minimizing delay jitter. Also, the efficiency of bandwidth allocation is improved by controlling the amount of bandwidth request depending on the queue length and packet arrival rate. We have analyzed the stability of the proposed mechanism based on a systematic approach. Using this analysis, we have derived a simple design guideline for the proposed algorithm and proved that the rate-based control in bandwidth request is essential for stability. Moreover, we have implemented the proposed algorithm in the OPNET simulator and have evaluated its performance in terms of optimality

of bandwidth allocation and ability of delay control. The simulation results confirm that the proposed algorithm strikes a balance between efficiency and QoS and that it provides a control knob for the delay by using the target delay. We expect that the proposed dual feedback mechanism can be a practical solution for the bandwidth request mechanism of real-time services and it can be widely extended to any centralized scheduling framework with dynamic bandwidth request.

REFERENCES

- [1] IEEE 802.16 WG, "IEEE standard for local and metropolitan area networks part 16: Air interface for fixed and mobile broadband wireless access systems, Amendment 2," *IEEE 802.16 Standard*, December 2005.
- [2] IEEE 802.16 WG, "IEEE standard for local and metropolitan area networks part 16: Air interface for fixed broadband wireless access systems," *IEEE 802.16 Standard*, June 2004.
- [3] K. Wongthavarawat and A. Ganz, "Packet scheduling for QoS support in IEEE 802.16 broadband wireless access systems," *International Journal of Communication systems*, vol. 16, pp. 81–96, 2003.
- [4] C. Cicconetti, L. Lenzini, E. Mingozzi, and C. Eklund, "Quality of service support in IEEE 802.16 networks," *IEEE Network*, vol. 20, pp. 50–55, March/April 2006.
- [5] G. Song, Y. Li, J. L. J. Cimini, and H. Zheng, "Joint channel-aware and queue-aware data scheduling in multiple shared wireless channels," in *Proceedings of IEEE Wireless Communications and Networking Conference (WCNC)*, vol. 3, pp. 1939–1944, 2004.
- [6] A. Sayenko, O. Alanen, J. Karhula, and T. Hämäläinen, "Ensuring the qos requirements in 802.16 scheduling," in *Proceedings of the 9th ACM international symposium on Modeling analysis and simulation of wireless and mobile systems (MSWiM)*, pp. 108–117, 2006.
- [7] Q. Liu, X. Wang, and G. B. Giannakis, "A cross-layer scheduling algorithm with QoS support in wireless networks," *IEEE Trans. on Vehicular Technology*, vol. 55, pp. 839–847, May 2006.
- [8] D. Niyato and E. Hossain, "Queue-aware uplink bandwidth allocation and rate control for polling service in IEEE 802.16 broadband wireless access networks," *IEEE Trans. on Mobile Computing*, vol. 5, no. 6, pp. 668–679, 2006.
- [9] OPNET WiMAX Model Development Consortium, "OPNET network simulator with WiMAX model." <http://www.opnet.com/WiMax>, 2006.
- [10] G. F. Franklin, J. D. Powell, and A. Emami-Naeini, *Feedback control of dynamic systems*. Addison-Wesley, 3rd ed., 1995.
- [11] N. Blaunstein, *Radio Propagation in Cellular Networks*. Artech House, 1999.
- [12] ITU-R Task Group 8/1, "Guidelines for evaluation of radio transmission technologies for IMT-2000," *Recommendation ITU-R M.1225*, 1999.
- [13] WiMAX Forum, "Mobile WiMAX - Part I: A technical overview and performance evaluation," *White Paper*, Aug 2006.