# Robust QoS Control for Single Carrier PMP Mode IEEE 802.16 System

Xiaofeng Bai, *Student Member, IEEE,* Abdallah Shami, *Member, IEEE,* and Yinghua Ye, *Member, IEEE*

*Abstract*—The IEEE 802.16 WirelessMAN standard provides a comprehensive quality-of-service (QoS) control structure to enable flow isolation and service differentiation over the common wireless network interface. By specifying a particular set of service parameters, the media access control (MAC) mechanisms defined in the standard are able to offer predefined QoS provisioning on per-connection basis. However, the design of efficient, flexible and yet robust MAC scheduling algorithms for such QoS provisioning still remains an open topic. This paper proposes a new QoS control scheme for single-carrier point-to-multipoint mode WirelessMAN systems, that enables the predefined service parameters to control the service provided to each uplink and downlink connection. By MAC-PHY cross-layer resource allocation, the proposed scheme is robust against particular wireless link degradation. Detailed simulation experiments are presented to study the performance and to validate the effectiveness of the proposed QoS control scheme.

*Index Terms*—Broadband Wireless Access Networks, Cross-layer Design, IEEE 802.16, Quality of Service, Scheduling.

## I. INTRODUCTION

**B**ROADBAND Wireless Access (BWA) technology is accepted as a promising solution for the next generation last-mile access systems [1], [2]. With less deployment time, lower maintenance cost and higher network scalability, this technology offers more revenue points and better market penetration than its wired counterpart. Supported by the robust link adaptation framework and flexible QoS control structure in the newly developed standard, BWA in particular provides a suitable solution for delivering broadband Internet services to foliage-populated rural areas, where wired infrastructure is economically infeasible, as well as to built-up urban areas, where existing wired access is bottlenecked.

The wireless metropolitan area network specifications, known as WirelessMAN, are defined in the IEEE standard 802.16-2004 [3]. This new release supersedes the previous IEEE standard 802.16-2001 [4] and its amendment [5], which documents the specifications for $2-11GHz$ band applications. The new standard provides a set of mechanisms to achieve reliable and link-adaptive high rate transmission over the wireless link. QoS is supported by the standard with a full set of parameters that permits service differentiation up to the connection level[1]. However, the implementation of QoS provisioning is left as vendor specific. In general, a viable

solution for such implementation should address the following issues:

1) How to inform the Base Station (BS) with the connection level bandwidth needs at each Subscriber Station (SS) in a timely manner.
2) How to allocate the limited radio resource among multiple SSs to ensure per-connection QoS delivery at each SS.
3) How to schedule the transmission among multiple connections sharing the same output channel to meet each connection's QoS requirement.

This study proposes a scheduling algorithm to address these fundamental issues by defining detailed operations performed at the BS and each SS. Moreover, the proposed algorithm is capable of providing robust connection level QoS control in various wireless channel conditions. Particularly, the contributions of this work can be summarized as below:

1) This work addresses the inherent conflict between per-connection bandwidth request and the in-band MAC signaling overhead defined in the standard. Namely, more connections in the system entail larger proportion of the limited radio resources to be used for signaling purpose. This severely threatens the individual service flows' QoS provisioning. Besides formulating the problem, this work proposes a MAC signaling approach that carries all the per-connection bandwidth requests with fixed overhead.
2) This work interprets the necessity of MAC-PHY cross-layer consideration for resource allocation in IEEE 802.16 systems. Many of the existing research work perform bandwidth allocation based exclusively on the MAC layer bandwidth requests. However, without considering the burst profiling applied over each wireless link, the arbitrary allocation dictated by MAC layer is not meaningful. Since each link is associated with an individual burst profile, one modulated symbol at the PHY layer provides a link-specific transmission capacity to the MAC layer. Therefore, the MAC-only based capacity allocation might not be fully accommodated by the PHY, where the number of symbols in one frame duration is fixed, or this capacity allocation might not be able to fill up all the symbols offered by the PHY.
3) This work proposes and verifies an efficient QoS control protocol design for the point-to-multipoint operation mode of single-carrier 802.16 system. The proposed scheme provides each connection differentiated service opportunities, such that the contracted QoS parame-

X. Bai and A. Shami are with the Department of Electrical and Computer Engineering, The University of Western Ontario, London, ON, Canada, N6A 5B9. Email: xbai6@uwo.ca, ashami@eng.uwo.ca.

Y. Ye is with Nokia Siemens Networks, Mountain View, CA, U.S.A. Email: yinghua.ye@nsn.com.

[1]The formal definition of *connection* is given in Section III-B. In this section it is not required to differentiate *connection* and *service flow*.

ters of each connection are constantly complied with. Moreover, by having the MAC layer resource allocation timely respond to the PHY layer link adaptation, this scheme is robust against wireless channel variation.

The rest of this paper is organized as follows: Section II briefly reviews some related research work. Section III presents an overview of the PHY and MAC layer specifications in the standard. Section IV depicts the main challenges for an efficient QoS control protocol design, as per the specifications included in the standard. Section V elaborates our proposed QoS control scheme. Simulation results are presented in Section VI and Section VII concludes this study.

## II. RELATED RESEARCH WORK

Currently, there are many research contributions on the design and analysis of IEEE 802.16 system as well as on cross-layer design of wireless communication systems. Among the IEEE 802.16 based research, particularly, the authors of [6] introduce a packet scheduling algorithm for QoS support in the 802.16 systems, where Fixed Allocation, Earliest Deadline First (EDF), Weighted Fair Queuing (WFQ) and Equal Sharing schemes are applied to the connections of different service types. However, this work considers only the bandwidth allocation to individual connections at the BS side. Since the bandwidth allocated to a SS is an aggregated grant for all connections at the SS, the ultimate QoS provisioning requires intelligent outbound transmission scheduling at each SS, which is not included in the paper. In [7], an uplink bandwidth allocation scheme for polling services is proposed and analyzed with advanced mathematical model and queuing analysis. Specifically, a Markov Modulated Poisson Process (MMPP) is applied to model the arrival process of polling services. Based on this queuing model, the proposed resource assignment scheme is evaluated. In contrast with [6], the limitation of [7] is that it only considers the outbound transmission scheduling at one SS assuming a fixed bandwidth allocation by the BS in every transmission frame, which eliminates the consideration for per-connection bandwidth request. The authors in [8] propose a more general QoS architecture, where both the bandwidth allocation at the BS and outbound transmission scheduling at the SS are outlined. Though many functions are introduced to provide QoS at the connection level, the overall architecture however does not address the stringent QoS parameters of individual connections. In [9], the authors propose a service flow management scheme that dynamically adjusts the bandwidth allocated to downlink and uplink subframes to achieve more flexibility and thereby to improve the system throughput. Moreover, in order to reduce frame creation work load at the beginning of each transmission frame, a special approach termed *Frame Registry Tree Scheduler* is proposed in [10]. In this work, the BS performs connection level bandwidth allocation "on-the-fly", when each connection's bandwidth request arrives at the BS. Therefore, it reduces the BS's offline computation time. Many aforementioned research works, e.g., [6], [8], [9], address service flow admission control at the BS. The common principle for admission control is to accept a new call only if: (a) active service flows' QoS requirements are not violated and, (b) the requesting service flow will be provided QoS guarantee. The QoS degradation model proposed in [11] slightly compromises this principle by degrading the maximum sustained traffic rate of some existing service flows to make room for the minimum reserved traffic rate of the requesting service flows.

Due to the radical variation of wireless channel, cross-layer protocol design for wireless communication system is currently under intensive study. Notable researches in this category include the cross-layer design framework proposed in [12], where the authors present this technique to maintain a certain level of packet loss rate and average throughput at the data link layer, by dynamically adjusting the target packet error rate at the PHY layer. Controlled by queuing status, the target packet error rate decides the Adaptive Modulation and Coding (AMC) mechanism at the PHY layer and thereby the desired data link layer performance is achieved. Moreover, [13] introduces a multi-layer design that involves App-MAC-PHY coordination. The goal of this design is to determine the optimum rate-adaptation within the MAC-PHY capacity region that maximizes the multimedia quality. Specifically, a video flow is divided into several sub-flows with each assigned a fraction of the total user rate. These sub-flows are given different rate distortion values thereby the transmitted video quality can be optimized according to the underlying rate adaptation mechanism in the MAC-PHY capacity region. Taking the IEEE 802.16 system as example, the authors of [14] noticed that due to the absence of link layer automatic repeat request (ARQ) for $10 - 66GHz$ applications in the standard, when the wireless channel condition degrades, the TCP layer delay increases while the data link layer delay decreases. This is because when retransmission becomes frequent and the TCP congestion window is reduced, the data link layer is given enough time to empty its queue. The authors therefore, concluded that the link adaptation scheme should be designed by considering the TCP layer throughput as well, instead of just measuring the PHY layer bit error rate. In this study, our proposed QoS control scheme performs resource allocation by taking into account both of MAC layer queuing and PHY layer burst profiling information. The proposed scheme is thus able to provide connection level QoS enforcement even when the wireless link condition degrades.

## III. OVERVIEW OF THE IEEE 802.16 WIRELESSMAN SYSTEM

The IEEE 802.16 WirelessMAN system contains one central base station and multiple subscriber stations in one architectural cell. The BS is responsible for communicating with each SS and regulating its behavior. Two operation modes are defined in the standard, i.e., point-to-multipoint (PMP) and mesh modes, along with different physical layer specifications. In this study, we focus on developing an efficient QoS control scheme for the $10 - 66GHz$ band Single-Carrier (SC) system in PMP operation mode.

### A. PHY specifications for WirelessMAN-SC system

The PHY layer operation is frame based and uses burst transmission format where the burst profiling for each SS is
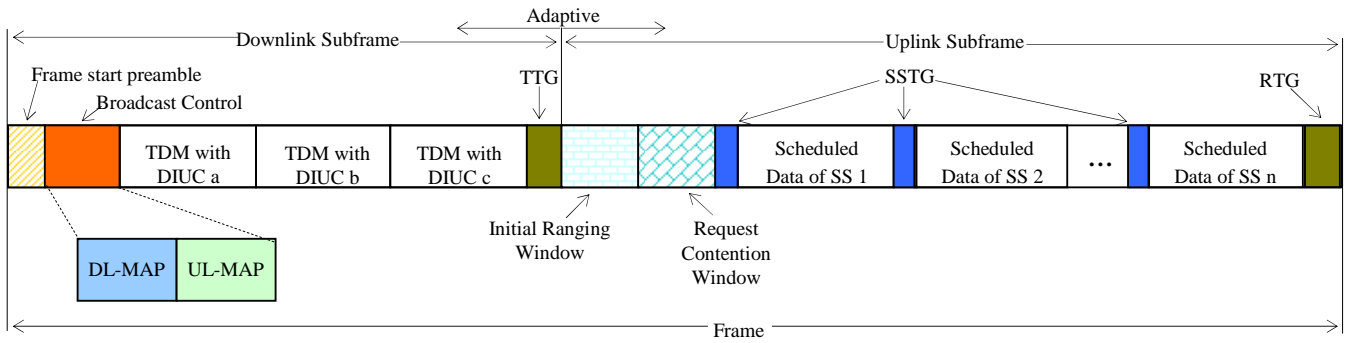
Fig. 1.   Illustration of TDD frame structure

adaptive and may change frame-by-frame.

For Time Division Duplex (TDD) based WirelessMAN systems, a transmission frame is defined as a fixed time duration that consists of two subframes, i.e., downlink and uplink, designated for BS-to-SS and SS-to-BS transmissions, respectively. As shown in Fig. 1, the downlink subframe begins with synchronization information and a frame control section. The frame control section is broadcast to all SSs and contains Downlink MAP (DL-MAP) and Uplink MAP (UL-MAP) messages that define the transmission burst profiles, including modulation and coding schemes as well as relevant timing information, for the following downlink and uplink transmissions, respectively. Following the frame control section is the downlink data destined to individual SSs. The downlink data is grouped into several transmission bursts using Time Division Multiplexing (TDM) technique. These transmission bursts are differentiated by their applied Downlink Interval Usage Code (DIUC), which represents a certain set of modulation and coding scheme for transmission. Downlink data bursts are broadcast by the BS to all SSs. The MAC in each SS listens to the downlink channel and looks for the MAC headers indicating data for this SS. The uplink subframe follows the downlink subframe. It may, but not necessarily in each frame, start with initial ranging and/or contention-based bandwidth request intervals, under the discretion of the frame scheduler at the BS. The initial ranging interval is used for new SSs to be registered into the system, while the bandwidth request contention interval is designed for connections carrying non-real-time applications to inform the BS of their bandwidth needs. Collisions in both contention-based intervals are resolved by binary exponential back-off algorithms [3]. The uplink data transmission is organized in Time Division Multiple Access (TDMA) fashion, where the uplink bursts are differentiated by the sending SSs. Each scheduled SS transmits into the uplink during its granted window using burst profile associated with the Uplink Interval Usage Code (UIUC) that is informed by the UL-MAP message in the frame control section.

### B. MAC scheduling services

While the PHY layer specifications are differentiated by the spectrum of usage, the standard is designed to evolve as a set of air interfaces based on a common MAC protocol. The operation of MAC services is connection-oriented. A connection is defined as a unidirectional mapping between BS and SS MAC peers for the purpose of transporting a service flow's traffic [3]. A service flow is a unidirectional flow of MAC Protocol Data Units (PDUs) with predefined QoS parameters. The QoS parameters defined for the service flow is therefore implicitly provided by the connection's unique Connection Identifier (CID). In order to accommodate applications with different service requirements, the standard defines four types of MAC scheduling service. Namely,

1) *Unsolicited Grant Service* (UGS)—designed to support real-time data streams consisting of fixed-size data packets issued at periodic intervals, such as T1/E1 and Voice over IP without silence suppression. The key QoS parameters relevant to this service type are reserved traffic rate, maximum latency and tolerated jitter.

2) *Real-time Polling Service* (rtPS)—designed to support real-time data streams consisting of variable-sized data packets that are issued at periodic intervals, such as Moving Pictures Experts Group (MPEG) video. The key QoS parameters relevant to this service type are minimum reserved traffic rate, maximum sustained traffic rate and maximum latency. The maximum latency here is only guaranteed within the scope of the minimum reserved traffic rate.

3) *Non-real-time Polling Service* (nrtPS)—designed to support delay-tolerant data streams consisting of variable-sized data packets for which a minimum data rate is required, such as FTP applications. The key QoS parameters relevant to this service type are minimum reserved traffic rate and maximum sustained traffic rate.

4) *Best Effort* (BE)—designed to support data streams for which no minimum service level is required and therefore may be handled on a space-available basis. The key QoS parameter relevant to this service type is only maximum sustained traffic rate.

Among these types of scheduling service, the bandwidth need for UGS service flows is time-invariant and delay sensitive. Therefore, bandwidth for UGS connections is offered by the BS periodically in a fixed pattern, and UGS connections do not request bandwidth from the BS. A specific study on the grant synchronization mismatch behavior of UGS service is conducted in [15], where the authors designed a *Grant Synchronization* approach to combat the unnecessary service delay derived from this synchronization mismatch. In this

study, we aim to focus on the QoS provisioning of other three polling-based service types, i.e., rtPS, nrtPS and BE, where the service opportunity is offered by the BS on *demand assignment* basis. In the following discussion, we do not include UGS service flows.

Each SS can request bandwidth for one of its connections by sending a stand-along MAC header (6 bytes) or by inserting a grant management subheader (2 bytes) into a MAC PDU, i.e., piggy-back request. For $10 - 66GHz$ system, Grant-per-SS (GPSS) is defined, where bandwidth is granted by the BS to a SS as an aggregate of grants in response to per-connection requests from the SS [3]. Therefore, the SS has the right to decide the ultimate allocation of the aggregated grant assigned by the BS. The per-connection bandwidth request allows the BS to allocate resources among SSs in a QoS-aware manner. For example, connections with stringent QoS requirements can be better serviced when the contention for radio resource arises, by appropriately granting the corresponding SSs. The per-SS granting mechanism, on the other hand, allows the SS to further adjust the usage of the granted bandwidth, hence to compensate the inherent drawback of polling operation due to information delay. Also the per-SS granting reduces signaling overhead in the downlink, without degrading system performance.

## IV. DESIGN CHALLENGES FOR EFFICIENT QOS CONTROL SCHEMES

It is clear that an efficient QoS control scheme should provide differentiated service to each connection, such that the various QoS parameters of individual connections are constantly complied with. Note that connections belonging to the same service type may have different QoS parameter settings. Although, the standard defines a variety of signaling mechanisms to facilitate efficient resource management and reliable QoS delivery, the detailed design of such a scheduling algorithm is left as vendor specific. It should be mentioned that in line with the standard, the downlink traffic can be also categorized into the above four scheduling service types for differentiated services considered by the BS. However, since the downlink connection information is locally available to the BS, the design challenges mainly lie in the scheduling of uplink connections. Taking the entire data control plane as a collaborative entity, the following aspects have to be considered for designing an efficient QoS control scheme:

1) The first challenge of this design is how to provide connection level QoS guarantee, assuming the per-connection bandwidth requests of each SS can be forwarded to the BS in a timely manner. By considering the scheduling service type and the real time bandwidth need of each connection in the network, the BS should appropriately partition the downlink and uplink subframes and grant every SS uplink transmission opportunity to satisfy the QoS requirements of connections running at the SS.

2) The second challenge of this design is how to provide the BS up-to-date information on the bandwidth need of each connection, i.e., how often should a connection send bandwidth request to the BS. With more up-to-date awareness of per-connection demand, the BS is able to manage the radio resource more efficiently. For example, delay sensitive connections are serviced faster and less bandwidth is wasted.

3) The third challenge of this design is how to reduce the signaling overhead used for per-connection bandwidth request. Due to the in-band signaling overhead for bandwidth request, as we have mentioned before, the QoS delivery has to be compromised when the number of connections in the network is large. Consider that the 16-bit CID defined in the standard permits up to $64K$ connections running simultaneously under the BS's administration, the bandwidth available for data transmission may decrease severely in peak hours. Prolonging the time interval for connections to send bandwidth request, i.e., reducing the request frequency, would apparently reduce the proportion of this overhead. However, this leads to suboptimal usage of the radio resource, which is the second challenge stated above.

In conclusion, the following principles should be taken into consideration when the new protocol is designed:

1) Guarantee solid compliance with each connection's QoS parameter settings, i.e., connections of service type:
   - rtPS—maximum latency, minimum reserved traffic rate and maximum sustained traffic rate;
   - nrtPS—minimum reserved traffic rate and maximum sustained traffic rate;
   - BE—maximum sustained traffic rate.

2) Optimize the freshness of BS's perception on each connection's bandwidth need.

3) Minimize operational overhead required for connections' bandwidth requests.

Bearing these points in mind, we propose a robust QoS control scheme that enables each connection's (both uplink and downlink) predefined QoS parameters to control the service provided to the connection. In the following discussion, we refer to this scheme as Single-Carrier Scheduling Algorithm (SCSA).

## V. ROBUST QOS CONTROL SCHEME FOR IEEE 802.16 WIRELESSMAN-SC SYSTEM

The proposed SCSA scheme operates via efficient collaboration of two functional blocks, i.e., the Uplink Request Management Agent (URMA) and the Frame Scheduling Unit (FSU), located at each SS and at the BS, respectively. In general, the functionality of each URMA is twofold:

1) communicating with the BS on each connection's bandwidth need at the SS, with minimal signaling overhead;

2) assisting to schedule uplink transmission at the SS, in QoS-aware manner.

Concurrently, the task of the FSU includes:

1) collecting each connection's bandwidth need information in the network;

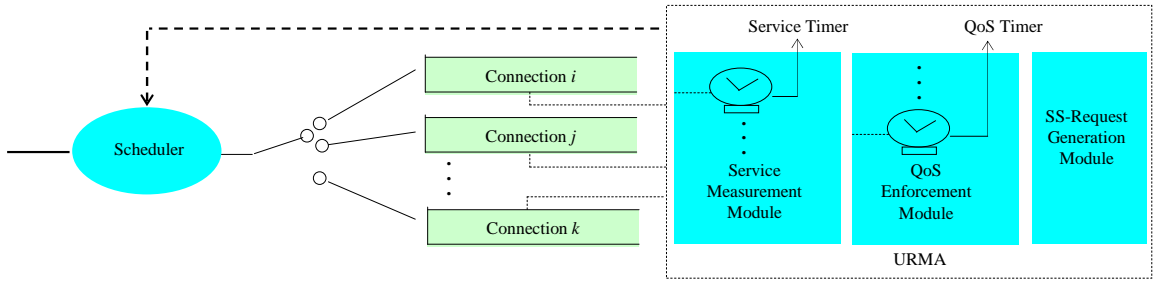2) performing per-SS resource allocation, based on real-time MAC-PHY cross-layer information;

Fig. 2. Illustration of URMA

3) defining a new frame, according to the resource allocation performed in 2);
4) assisting to schedule downlink transmission at the BS, to provide QoS for downlink connections.

The basic idea of this SCSA scheme for solving the third challenge, i.e., minimizing the operational overhead required for bandwidth requests, is to move some connection level functionalities performed by the BS to each SS, thereby to reduce the signaling overhead. Specifically, an URMA is installed at each SS and collaborates with the SS's local scheduler. This URMA partially processes each connection's bandwidth request at the SS and only sends the packed information that is necessary for the BS to reach an optimal solution for its per-SS based resource allocation. Upon receiving grant from the BS, the URMA finalizes the resource allocation with respect to each connection, such that the QoS requirements of each connection are satisfied. In this configuration, the BS operates as a "server" while each SS acts like a "workstation", with the URMA emulating a "Java applet" in between. After preprocessing the bandwidth request of each connection running at the SS, the URMA generates up to three per-SS bandwidth requests that will be sent to the BS. These requests are labeled with different priority values and are sent at the end of the uplink transmission window assigned to the SS. In this way, the overhead required for bandwidth request is limited to be only SS-relevant and independent from the number of connections running at the SS. Therefore, these per-SS bandwidth requests can be generated and sent to the BS in every transmission frame, in order to feed the BS with the most recent information on the bandwidth need of each connection, i.e., to optimize the freshness of BS's perception on each connection's bandwidth need. Hence, the second principle is complied with. Next we will focus on the first and most challenging principle, i.e., to provide connection level QoS guarantee.

### A. Uplink Request Management Agent

An URMA consists of three modules, i.e., Service Measurement module, QoS Enforcement module and SS Request Generation module, as shown in Fig. 2. The functionalities performed by each module are described below.

*1) Service Measurement module:* This module computes the instant bandwidth request of each connection at the end of each uplink transmission window of the SS, according to the connection's queue length and MAC headers required

to transmit the backlogged traffic. This bandwidth request is upper-bounded by the connection's, say connection $i$'s, eligible bandwidth request denoted as $r_i^e$. Namely, if the required bandwidth for servicing the backlogged traffic is less than $r_i^e$, the connection can request its required amount. Otherwise, it is only eligible to request bandwidth of amount $r_i^e$. The eligible bandwidth request of connection $i$, i.e., $r_i^e$, is computed as follows:

$$r_i^e = \max\left\{ \frac{R_i^{max}}{8} \times [t - S_i(t)],\, 0 \right\} \quad (1)$$

where $R_i^{max}$ is the maximum sustained traffic rate (in bits/second) of this connection introduced in Section III-B, and $t$ is the SS's system time. The service measurement module maintains a service timer for each non-UGS connection running at the SS, and $S_i(t)$ in (1) is the value of connection $i$'s service timer at time $t$. The service timer applies the concept of *Virtual Time* proposed by L. Zhang in [16]. This service timer is synchronized with the system clock when a connection is established and ticks with the following value upon the service of each PDU in the corresponding connection:

$$A_i = \frac{8 \times B_i}{\rho_i} \quad (2)$$

where $A_i$ is the increment of connection $i$'s service timer, upon the service of a PDU with size of $B_i$ bytes. $\rho_i$ is the measurement rate (in bit/second) for connection $i$. For the service timer, the value of $\rho_i$ is $R_i^{max}$. Instead of stamping each arriving packet with its virtual time as introduced in [16], here we use a separate variable to record a connection's virtual time. The reason, as will be shown later, is that in our case some packets may be dropped-front due to violation of their predefined maximum latency. These packets do not consume bandwidth and thus should not consume virtue time. The service timer is applied to threshold a connection's service rate within its maximum sustained traffic rate. Namely, when the service timer meets the SS's system time, the corresponding connection is not *eligible* to receive any service (thus to request any bandwidth) at current time, as interpreted in (1).

*2) QoS Enforcement module:* This module maintains a QoS timer for each rtPS and nrtPS connection running at the SS. The QoS timer is also initialized with the SS's system clock, however ticks with the value decided by different measurement rate defined in (2). For the QoS timer, the value of $\rho_i$ should be $R_i^{min}$, i.e., the minimum reserved traffic rate (in bits/second) of connection $i$. The reason for designing two virtual time

values for a connection requiring minimum service rate (i.e., rtPS or nrtPS connection) is to have the received service rate "guaranteed but not limited" by its minimum reserved traffic rate. This will be further clarified later by the discussion of outbound transmission scheduling. The QoS timer enforces the connection's service rate to meet a guaranteed value. This is achieved by differentiating the bandwidth guaranteed part and non-bandwidth guaranteed part of the connection's bandwidth request, as explained below.

The operations performed by the QoS enforcement module include two steps:

1) For each rtPS or nrtPS connection $i$, the QoS enforcement module divides its bandwidth request $r_i$ into bandwidth guaranteed (BG) part and non-bandwidth guaranteed (NBG) part, i.e., $r_i^{BG}$ and $r_i^{NBG}$, according to the corresponding value of QoS timer, denoted as $Q_i(t)$, i.e.:

$$r_i^{BG} = \begin{cases} \min\left\{r_i, \frac{R_i^{min}}{8} \times [t - Q_i(t)]\right\} & Q_i(t) < t \\ \\ 0 & Q_i(t) \geq t \end{cases} \tag{3}$$

$$r_i^{NBG} = r_i - r_i^{BG} \tag{4}$$

where $r_i$ is connection $i$'s bandwidth request (in bytes) computed by the service measurement module, i.e., the minimum value between connection $i$'s required bandwidth and its eligible bandwidth request $r_i^e$. The bandwidth request of a BE connection is always NBG, as there is no service guarantee for this connection.

2) For each rtPS connection $i$, the QoS enforcement module further divides its $r_i^{BG}$ into imminent part $r_i^{BG-im}$ and non-imminent part $r_i^{BG-nim}$, according to the contracted maximum latency of this connection. The operation of this step is based on the following facts: if the maximum latency deadline of a packet will appear in frame $n+2$, the latest frame to transmit this packet is frame $n+1$, in order to guarantee its maximum latency. Therefore, the bandwidth needed to transmit this packet has to be requested in frame $n$. The $r_i^{BG-im}$ part of $r_i^{BG}$ computed in frame $n$ is thus the bandwidth needed for transmitting packets whose maximum latency deadline appears before the end of frame $n+2$. The $r_i^{BG-nim}$ part of $r_i^{BG}$ is simply $r_i^{BG} - r_i^{BG-im}$. Fig. 3 illustrates the bandwidth request of connections processed by the QoS enforcement module.

*3) SS-Request Generation module:* This module generates up to three per-SS bandwidth requests, depending on the service type of connections running at the SS and the output of the QoS enforcement module. The per-SS bandwidth requests are prioritized, in order to enable service differentiation at the

BS. These three bandwidth requests are defined as:

$$r_{SS}^{P0} = \sum_{i \in M} r_i^{BG-im} \tag{5}$$

$$r_{SS}^{P1} = \sum_{i \in M} r_i^{BG-nim} + \sum_{j \in N} r_j^{BG} \tag{6}$$

$$r_{SS}^{P2} = \sum_{i \in M} r_i^{NBG} + \sum_{j \in N} r_j^{NBG} + \sum_{k \in L} r_k \tag{7}$$

where $r_{SS}^{Pi}$ ($i = 0, 1, 2$) denotes the per-SS bandwidth request of priority level $i$, $M$, $N$ and $L$ are the sets of rtPS, nrtPS and BE connections running at the SS, respectively.

Before going further, it is helpful to clarify the background basis whereby these three prioritized requests are defined:

1) From the QoS viewpoint, the $r^{BG}$ part of a rtPS or nrtPS connection is prioritized over its $r^{NBG}$ part and the latter is equivalent to the bandwidth request of a BE connection;

2) From the QoS viewpoint, the $r^{BG-im}$ part of a rtPS connection is prioritized over its $r^{BG-nim}$ part and the latter is equivalent to the $r^{BG}$ part of a nrtPS connection.

### B. Frame Scheduling Unit

Upon receiving the prioritized bandwidth requests from SSs in the uplink subframe, a new frame is generated by the FSU at the BS. This requires the involvement of three functional modules, i.e., the Downlink Request Management (DRM) module, the Resource Allocation module and the Frame Creation module. The operations performed in these modules are described below.

*1) Downlink Request Management module:* The DRM module functions similarly to the URMA at a SS, except that the prioritized bandwidth requests for downlink are associated with each downlink burst. Namely, downlink connections using the same DIUC in the next frame are grouped together by a common set of prioritized bandwidth requests. Note that uplink connections only in the same SS are grouped together by a common set of prioritized bandwidth requests. In the following discussion we refer to connections sharing a common set of prioritized bandwidth requests as a Scheduling Group (SG).
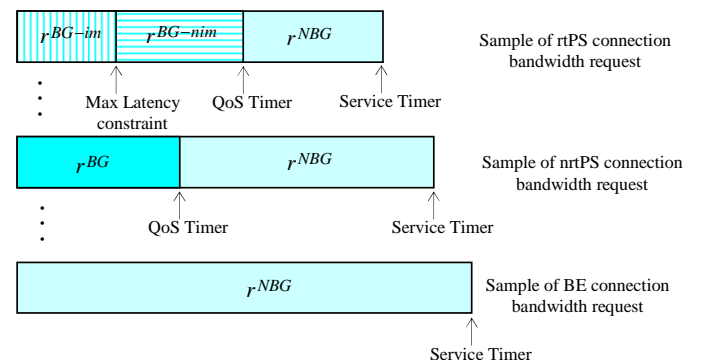


Fig. 3. Illustration of connections' bandwidth request processed by the QoS enforcement module

*2) Resource Allocation module:* This module allocates transmission capacity to each SG, according to their prioritized bandwidth requests. In practical wirelessMAN systems the wireless channel is time-variant, in the sense that the PHY layer performance over a wireless link may vary severely with time, if the applied burst profile is fixed. For this reason, the standard specifies adaptive burst profiling mechanism at the PHY layer, to enable diverse vendor-level link adaptation design. Namely, when the received *signal-to-noise ratio* (SNR) over particular wireless link is intolerably attenuated, the MAC layer data rate can be temporarily compromised-down for better PHY layer transmission performance, by switching the applied burst profile to a more robust but less efficient version, i.e., bursts carrying less MAC layer bits with one PHY layer symbol. Implementing this adaptive burst profiling mechanism imposes mandatary modification of the MAC layer protocols, as this implies that each SG may place very different constraints on resources needed, at particular time instant, for transmitting the same amount of MAC layer data. For example, SGs applying QPSK modulation require four modulated symbols while SGs applying 16-QAM modulation utilize two modulated symbols, to carry 4-bits data in the MAC PDU (assuming the same code rate of $1/2$ ). Therefore, the prioritized bandwidth requests of each SG have to be converted into symbol needs according to individual burst profile being used and hence the resource allocation is finally accomplished by symbol assignment. Simply put, let $\eta_i$ be the combined modulation-coding symbol efficiency, i.e., number of bits in the MAC PDU that can be carried by one modulated symbol, of the $i^{th}$ SG. For example, $\eta = 1\,bit/symbol$ for SGs transmitting with QPSK modulation and $1/2$ code rate, while $\eta = 3\,bits/symbol$ for SGs transmitting with 16-QAM modulation and $3/4$ code rate. Then the number of symbols needed to accommodate the $Pi\,(i = 0,\,1,\,2)$ request of the $j^{th}$ SG is given as:

$$rs_j^{Pi} = \frac{8 \times r_j^{Pi}}{\eta_j} \tag{8}$$

where $r_j^{Pi}$ and $rs_j^{Pi}$ represent the $Pi$ request in bytes of the $j^{th}$ SG and the converted symbol need (may be a fractional value) of this request, respectively.

The symbol assignment follows strict priority rule and is defined as follows:

1) Symbol need of $P0$ request is considered first, then $P1$ followed by $P2$ requests, according to the number of symbols available for data transmission in a frame duration. Here the symbol need of $Pi\,(i = 0,\,1,\,2)$ requests includes the symbol need of every SG's $Pi$ request (both downlink and uplink).

2) During this symbol assignment operation, if the remaining number of symbols can be assigned by the BS is enough to satisfy the symbol need of all $Pi$ requests, these requests are fully accommodated.

3) If the remaining number of symbols is inadequate to satisfy the symbol need of all $Pi$ requests, every SG submitting a $Pi$ request shares the available symbols in

proportion to its symbol need, i.e.:

$$g_j^{Pi} = \frac{rs_j^{Pi}}{\displaystyle\sum_{k=1}^{n} rs_k^{Pi}} \times N_s \quad \left(N_s < \sum_{k=1}^{n} rs_k^{Pi}\right) \tag{9}$$

where $g_j^{Pi}$ is the symbols "earned" by the $Pi\,(i = 0,\,1,\,2)$ request of the $j^{th}$ SG, $N_s$ is the remaining number of symbols can be assigned by the BS, and $n$ is the number of SGs submitting $Pi$ requests.

4) If the symbol needs ($P0$, $P1$, and $P2$) of every SG have been fully accommodated, the remaining symbols are assigned to each SG, in proportion to the number of connections included in the SG.

5) The total number of symbols assigned to a SG (may be a fractional value) is the sum of symbols obtained by its $P0$, $P1$, $P2$ requests and any symbols assigned in step 4), if applicable. This number has then to be adjusted into integer values as feasible output of the resource allocation module. Particularly, the integer number of symbols finally assigned to the $j^{th}$ SG, i.e., $a_j$, is determined as:

$$a_j = \left\lfloor \sum_{i=0}^{2} g_j^{Pi} + g_j^{ex} + f_{j-1} \right\rfloor \tag{10}$$

$$f_j = \sum_{i=0}^{2} g_j^{Pi} + g_j^{ex} + f_{j-1} - a_j \tag{11}$$

where $g_j^{ex}$ is the extra symbols possibly assigned to this SG in step 4). The operator $\lfloor \cdot \rfloor$ in (10) takes the truncated integer value of the operand. $f_j$ is an adjustment factor that redistributes fractional symbol assignments among SGs.

*3) Frame Creation module:* The frame creation module converts the above symbol assignment result into timing information in terms of Physical Slot (PS) or minislot [3], and then a new frame is created. The slotization process performed in this module (one PS is four symbols long) involves similar truncation and fractional slot adjustment approach to the one described by (10) and (11).

## C. Outbound transmission scheduling

Within the granted transmission window in the uplink sub-frame, the scheduler at the SS selects packets for transmission, according to each connection's bandwidth request and the SS's prioritized bandwidth requests stored in the URMA by the last frame. Specifically, the $P0$ request is honored first, then $P1$ followed by $P2$ requests, if there is more bandwidth available. Based on each connection's bandwidth request of the last frame stored in the URMA, when the scheduler has to select the next packet for transmission among multiple connections, some criterion should be defined. Depending on the type of prioritized request to which the bandwidth is honored, three different criteria are applied:

1) When multiple connections contend for bandwidth honored to the $P0$ request, the packet with the most imminent maximum latency deadline is selected. In particular,

the scheduler compares the allowed delay of each *head-of-line* (HOL) packet in these connection queues and selects packet with the minimum allowed delay. However, any packet with expired maximum latency deadline (e.g., in overloaded connections) will be dropped at the front of the connection queue.

2) When multiple connections contend for bandwidth honored to the $P1$ request, the HOL packet of connection with the earliest QoS timer is selected. This criterion ensures the best fairness among contending connections, as it forces the QoS timer of these connections to keep up with others.

3) Bandwidth honored to the $P2$ request is equally shared by each connection in a round-robin fashion, despite of the service type and requested bandwidth of the connections. This criterion continues till the granted transmission window expires, i.e., even there may be no request to be honored. However, any connection whose service timer $S_i(t)$ reaches the SS's system time $t$, i.e., the corresponding maximum sustained traffic rate is met, will be neglected in the current round.

With the service of each PDU in a connection, the service timer ticks accordingly. However, only consuming bandwidth honored to the $P0$ and $P1$ requests incurs the advancement of corresponding QoS timer. This enables the received service rate of a rtPS or nrtPS connection to be "guaranteed but not limited" by its minimum reserved traffic rate.

The scheduler at the BS performs similarly to the scheduler at the SS for downlink outbound transmission scheduling, except that the operation scope is a downlink burst window. Namely, connections belonging to the same downlink SG are all considered when the next packet is selected for transmission.

## VI. SIMULATION EXPERIMENTS

In this section, we verify the effectiveness and robustness of the proposed SCSA scheme using simulation experiments. In the simulation, we apply the following assumptions:

1) Each SS is equally distant from the BS.
2) Line-of-sight is available over each wireless link.
3) The variation of any wireless channel is independent from others.
4) All connections have been admitted and the system is not overloaded, i.e., the sum of minimum reserved traffic rate for existing rtPS and nrtPS connections does not exceed the system capacity[2].
5) There is no service flow arrival and departure occurred throughout the simulation.

Specifically, we will test the system performance in two scenarios:

1) Regular operation—in this scenario, we compare the service parameters measured in the simulation with the corresponding QoS parameters predefined for each connection. This experiment aims to test the connection

[2]Some connection admission control schemes in literature could be applied here for this purpose, e.g., [11].

level QoS enforcement capability of the new SCSA scheme. Particularly, we are expecting that: a) The average throughput of any rtPS, nrtPS or BE connection is no more than its maximum sustained traffic rate; b) The average throughput of any rtPS or nrtPS connection is no less than its minimum reserved traffic rate, or equal to the offered rate of this connection; c) PDUs constituting the minimum reserved traffic rate of a rtPS connection are delayed less than the maximum latency constraint of this connection. Namely, if a rtPS connection is offered larger traffic intensity than its minimum reserved traffic rate, the proportion of PDUs violating latency constraint is no more than the overloading traffic portion above the minimum reserved traffic rate.

2) Link degradation—in this scenario, we simulate the situation where the wireless link condition from particular SS to the BS degrades, which invokes more reliable but inefficient burst profile applied over this link. This experiment is designed to purposely test the robustness of the MAC-PHY cross-layer resource allocation design in the new SCSA scheme. To further challenge this robustness feature of the new scheme, we also increase the minimum reserved traffic rate parameter of some connections at the SS where the link condition degrades, while keeping traffic offering the same as in scenario one. Particularly, we are expecting that: a) QoS parameters of each connection are still well maintained as depicted in scenario one, even for connections requiring augmented service guarantee over the deteriorated channel; b) The capacity loss due to particular wireless link degradation is averaged over the entire network, such that the connection level service variation is minimized.

### A. Uncontrolled scheduling algorithm (UCSA)

In order to distinguish the advantages of the new SCSA design, as comparison, we also simulated another scheduling algorithm that we refer to as *uncontrolled scheduling algorithm*. This UCSA scheme merely implements the standard specifications, without introducing vendor-level QoS control intelligence. The key properties of this UCSA include:

1) Each uplink connection sends individual bandwidth request over the uplink transmission window assigned to the parent SS, to update the BS's perception on its bandwidth need.

2) The BS's MAC converts symbol need of each connection, according to its request, without knowing any burst profile change at the PHY layer. Namely, the symbol needs are estimated only based on the link capacity information when each connection is established.

3) The symbol assignment is performed according to strict priority rule as detailed in V-B2. Instead of $P0$, $P1$ and $P2$ requests as in SCSA, here the BS considers individual symbol need of each connection, i.e., symbol needs of all rtPS connections are accommodated first, then all nrtPS connections followed by all BE connections, if more symbols are available.

TABLE I
CIDS AND QoS PARAMETER SETTINGS FOR THE SIMULATION

| | CID | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | rtPS | | | nrtPS | | | BE | | |
| | Downlink | Uplink | | Downlink | Uplink | | Downlink | Uplink | |
| SS1 | 1 | 4 | 5 | 2 | 6 | 7 | 3 | 8 | 9 |
| SS2 | 10 | 13 | 14 | 11 | 15 | 16 | 12 | 17 | 18 |
| SS3 | 19 | 22 | 23 | 20 | 24 | 25 | 21 | 26 | 27 |
| SS4 | 28 | 31 | 32 | 29 | 33 | 34 | 30 | 35 | 36 |
| SS5 | 37 | 40 | 41 | 38 | 42 | 43 | 39 | 44 | 45 |
| SS6 | 46 | 49 | 50 | 47 | 51 | 52 | 48 | 53 | 54 |
| SS7 | 55 | 58 | 59 | 56 | 60 | 61 | 57 | 62 | 63 |
| SS8 | 64 | 67 | 68 | 65 | 69 | 70 | 66 | 71 | 72 |
| SS9 | 73 | 76 | 77 | 74 | 78 | 79 | 75 | 80 | 81 |
| SS10 | 82 | 85 | 86 | 83 | 87 | 88 | 84 | 89 | 90 |
| Offered rate ($Mbps$) | 0.8 | 0.8 | 1.2 | 0.8 | 0.8 | 1.2 | 1.2 | 1.2 | 1.2 |
| Max sustained rate ($Mbps$) | 1.0 | | | 1.0 | | | 1.0 | | |
| Min reserved rate ($Mbps$) | 0.5 | | | 0.5 | | | N/A | | |
| Max latency ($ms$) | 5 | | | N/A | | | N/A | | |

4) When the SS or BS schedules for uplink or downlink outbound transmission, the scheduler follows class level strict priority rule. Namely, any nrtPS connection could be serviced only when every rtPS connection queue is evacuated, and no BE connection could be serviced if any rtPS or nrtPS packet is backlogged. During servicing multiple connections of the same service type (rtPS, nrtPS or BE), the scheduler always selects the connection queue with the most severely delayed HOL packet for transmission. This policy favors overloaded connections without quantitative control by their contracted QoS parameters.

*B. Simulation model*

We developed a simulation environment by *ns-2* [17]. The simulated network consists of one BS and ten SS nodes (numbered from 1 to 10) locating at $2.5km$ away from the BS. In the downlink direction there are ten rtPS, ten nrtPS and ten BE connections originated from the BS; while in the uplink direction there are two rtPS, two nrtPS and two BE connections originated from each SS. The CIDs and QoS parameter settings of each connection are listed in Table I. In the simulation we apply three burst profiles to simulate the transmission with different reliability and efficiency. Namely, QPSK modulation with $1/2$ code rate, 16-QAM modulation with $3/4$ code rate and 64-QAM modulation with $2/3$ code rate, which are shortened in the following discussion as QPSK, 16-QAM and 64-QAM, respectively. We start the simulation with 16-QAM applied over each wireless link, except links from SS6 and SS8 to the BS, where 64-QAM is applied, to simulate the case where different channel conditions coexist in the network. The frame control section is always transmitted by QPSK for the best reliability, throughout the simulation.

The frame duration and system symbol rate for the simulation are set to $1ms$ and $20MBaud$, respectively. The standard specifies that no link layer ARQ should be applied in the $10-66GHz$ SC system, therefore we do not involve ARQ error control in the simulation. Following the standard, the frame start preamble and TDMA preamble preceding each uplink burst are 32-symbol and 16-symbol long, respectively. The

*transmit/receive transition gap* (TTG) and *receive/transmit transition gap* (RTG) are both set to $20\mu s$, which is long enough for the signal to propagate between the source and destination nodes. We neglect the BS processing time for generating the new frame. The PDU size is set to 70 bytes (minimum Ethernet packet size plus 6 bytes MAC header) and each connection queue has a limited buffer size for 50 PDUs. The *ns-2* built-in Exponential traffic model is applied to simulate the traffic flow offered to each connection. Each simulation experiment runs for 10 seconds, i.e., $10^4$ frames long, and the following results are observed.

*C. Simulation results and discussions*

*1) Regular operation:* In this scenario we run the simulation as described above to test, by both UCSA and SCSA, the compliance of measured service parameters for each connection with its predefined QoS parameters. For different scheduling service types, the performances of three sample connections in different SSs are evaluated. Particularly, connection 40 in SS5 for rtPS type, connection 52 in SS6 for nrtPS type, and connection 45 in SS5 for BE type are selected for illustration. Fig. 4(a)-4(c) visualize the average throughput and offered traffic intensity of these connections. The curves in the figure (and following figures) present the up-to-date averaged throughput and offered traffic intensity of the sampled connection. Therefore, with the evolvement of time, each curve becomes flatter and more precisely reveals the overall offered/serviced traffic rate of the connection. It is explicitly shown in the figure that:

1) With UCSA, rtPS connection 40 is suffced with enough bandwidth offering as in Fig. 4(a). While in Fig. 4(b), nrtPS connection 52 is provisioned higher throughput than its maximum sustained traffic rate. Since outbound transmission scheduling by UCSA follows class level strict priority rule and the bandwidth demand of nrtPS connections is not fully accommodated, e.g., connection 52, all BE connections in the network are constantly starved, as illustrated by connection 45 in Fig. 4(c). On the contrary, with SCSA, the average throughput of

(a) rtPS connection 40 throughput



(b) nrtPS connection 52 throughput



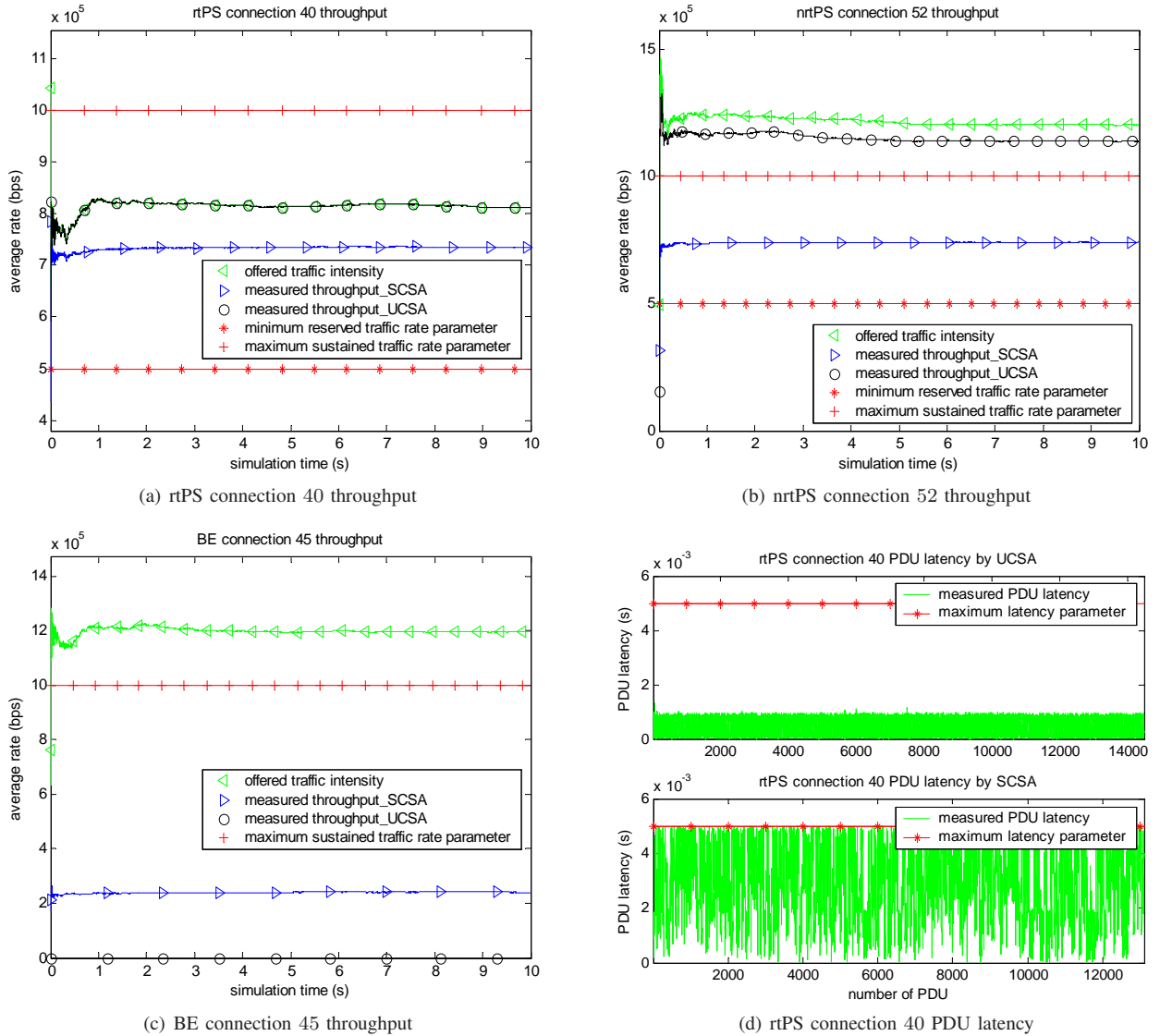(c) BE connection 45 throughput



(d) rtPS connection 40 PDU latency

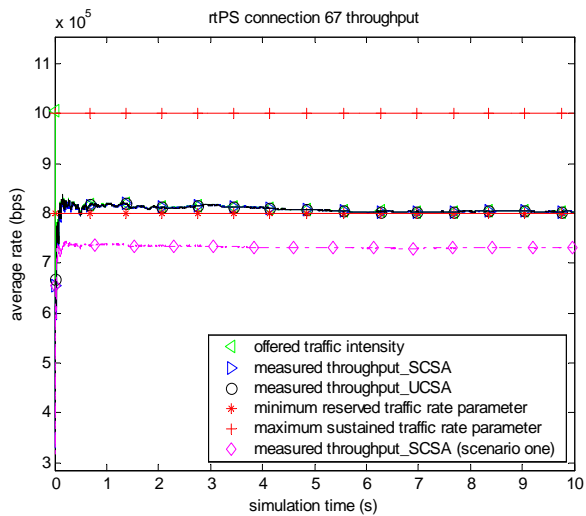Fig. 4.   Simulation results (scenario one)

each sampled connection is always constrained below the maximum sustained traffic rate of the connection, yet with no bandwidth starvation on any connection. The reason is that, by performing appropriate packet level admission control, SCSA does not improvidently favor overloaded rtPS or nrtPS connections and bias BE connections.

2) Without indulging overloaded connections, SCSA is still capable of ensuring the minimum reserved traffic rate of each rtPS or nrtPS connection. It is revealed in Fig. 4(a) and Fig. 4(b) that, with SCSA, the throughput curve of connection 40 and 52 are well maintained above, but not limited at, the minimum reserved traffic rate of each connection.

Moreover, in the simulation we have observed that with SCSA the PDU dropping probability of connection 40 due to latency violation is $9.54\%$. This value is less than the overloading traffic portion of this connection, i.e., $(0.8 - 0.5)/0.8 = 37.5\%$ (see Table I), which indicates that the proportion

of PDUs violating latency constraint is no more than the overloading traffic portion above the minimum reserved traffic rate. Therefore, the expected performance of the new SCSA design in this simulation scenario has been verified.
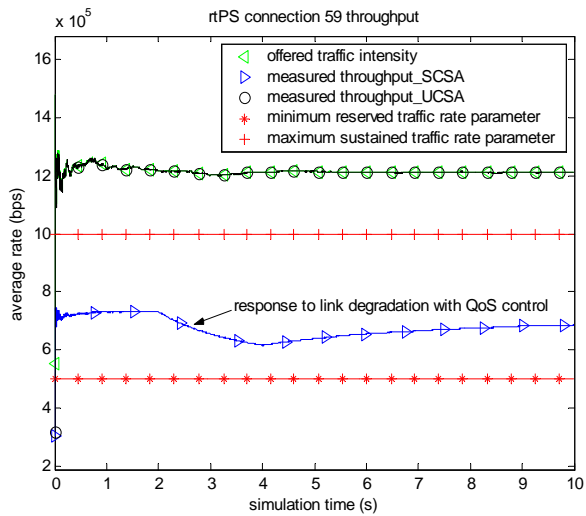
As complementary illustration, Fig. 4(d) shows the PDU latency of rtPS connection 40 obtained by both UCSA and SCSA. Since UCSA unconditionally prioritizes rtPS connections for outbound transmission scheduling, most PDUs in connection 40 are forwarded into the channel within one frame duration $(1ms)$ after arrival. With SCSA, however, some overloading PDUs cannot be serviced promptly by service opportunities offered to the guaranteed traffic portion of this connection. These PDUs will take chances to share the non-guaranteed service opportunities offered to the SS, with other nrtPS and BE connections in the SS. Therefore, some PDUs may experience large service latency that is close or equal to the maximum latency constraint of the connection, as seen in the bottom figure. It should be noted that UCSA offers less latency to rtPS connections with the cost of inappropriately
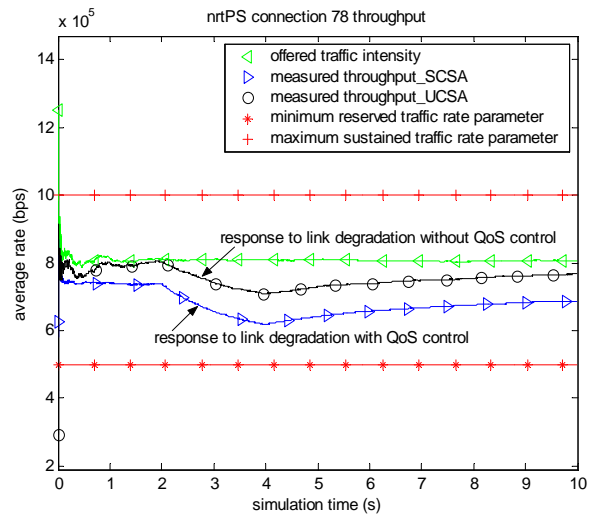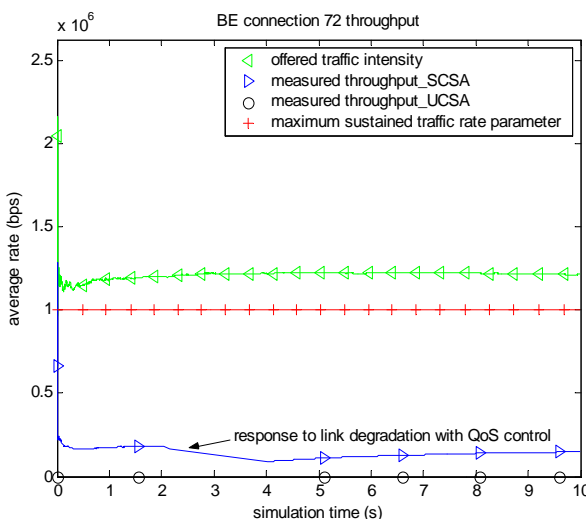
(a) rtPS connection 67 throughput
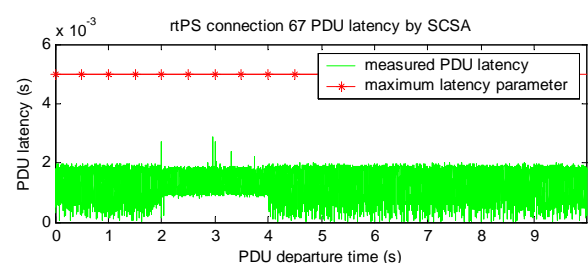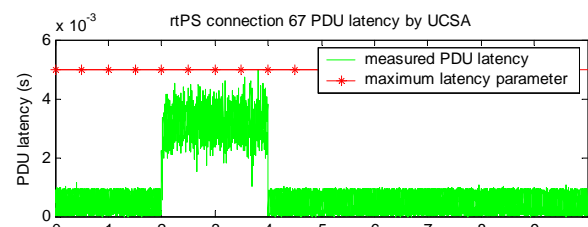


(b) nrtPS connection 69 throughput



(c) rtPS connection 59 throughput



(d) nrtPS connection 78 throughput



(e) BE connection 72 throughput



(f) rtPS connection 67 PDU latency

Fig. 5. Simulation results (scenario two)

starving BE connections, while SCSA prevents this misbehavior without violating the QoS parameters of any connection.

*2) Link degradation:* In this scenario we evaluate the robust cross-layer resource allocation design of SCSA against wireless link degradation and service augmentation. Specifically, at time 2.0 second the wireless link from SS8 to the BS degrades and the PHY layer link adaptation mechanism is invoked to change the burst profile applied over this link from 64-QAM to QPSK. Then at time 4.0 second the link recovers and QPSK is converted back to 64-QAM. We are interested in the response of different connections to the capacity loss caused by this link degradation. Moreover, we increase the minimum reserved traffic rate parameter of rtPS and nrtPS connections in SS8 from $0.5Mbps$ to $0.8Mbps$, while keeping other simulation settings the same as in scenario one (see Table I). We hope to see that even during link degradation, the augmented service guarantee of these connections is still provisioned by SCSA. Fig. 5(a) and 5(b) show the measured throughput of rtPS connection 67 and nrtPS connection 69 in SS8, i.e., where the link degradation occurs. Fig. 5(c)-5(e) show the service provisioning of rtPS connection 59 in SS7, nrtPS connection 78 in SS9 and BE connection 72 in SS8, which are randomly selected to represent connections of different service type locating in and out of SS8. The simulation reveals that:

1) With SCSA, the service curve of each rtPS or nrtPS connection is well maintained at/above its minimum reserved traffic rate, and the maximum sustained traffic rate parameter still upper-bounds service provisioned to any connection (rtPS, nrtPS or BE), as observed in scenario one. To illustrate service augmentation, here the measured throughputs of connections 67 and 69 in scenario one are also plotted for comparison. It is shown that the throughputs of these two connections are both increased from below $0.8Mbps$ in scenario one to $0.8Mbps$ in scenario two, irrespective of the link degradation. In contrast, UCSA is unable to offer guaranteed service to connection 69 (Fig. 5(b)) due to link degradation at SS8. Simultaneously, however, connection 59 (Fig. 5(c)) is receiving extra service beyond its maximum sustained traffic rate constraint, and connection 72 (Fig. 5(e)) is constantly starved as in scenario one.

2) The service curves offered by SCSA in Fig. 5(c), 5(d) and 5(e) distinguish a clear service drop, but without QoS violation, starting at 2.0 second and terminating at 4.0 second. This service drop is contributed to rtPS and nrtPS connections in SS8 for maintaining their augmented and guaranteed throughputs, as seen in Fig. 5(a) and 5(b). Hence we can conclude that SCSA averages the risk of particular link degradation over the entire network, by compromising service provisioned to the non-guaranteed traffic portion of each connection. This robustness feature greatly reduces the probability of QoS violation for connections involving in the link degradation.

3) It is noticed that with UCSA, the service curve of connection 78 (Fig. 5(d)), as response to the serious bandwidth starvation of connection 69 (Fig. 5(b)),

also slightly drops during the link degradation. This is resulted from the anarchical bandwidth contention occurred in UCSA. Particularly, when the link degrades, SS8's PHY layer is unable to serve the amount of traffic intended by the BS's MAC, using the assigned transmission window. This enforces worse traffic backlog at each connection queue at SS8[3]. By presenting larger bandwidth requests to the BS's MAC, the increasing traffic backlog at SS8 eventually diminishes the resource allocation to other SSs and hence the service offering to their housed connections, such as connection 78 at SS9 mentioned above. It should be noted that during link degradation, the service drop of connection 78, is a passive response to the QoS violation of connection 69 by UCSA, while is a voluntary contribution to the QoS compliance of connection 69 by SCSA.

Fig. 5(f) shows the PDU latency against PDU departure time of rtPS connection 67, obtained by both UCSA and SCSA. As interpreted in scenario one, arriving packets are serviced by UCSA within one frame duration, before and after the link capacity loss at SS8. Since some arriving packets have to be backlogged in the queue during the link capacity loss, PDU latencies are appreciably raised up to between two and four frame durations, as shown in the top figure. Differentiated from UCSA, SCSA ensures the minimum reserved traffic portion of the connection to be serviced within two frame duration, i.e., one for polling and one for service. Enabled by MAC-PHY cross-layer resource allocation, SCSA maintains this commitment even during the link capacity loss, as seen in the bottom figure. The immaterial effects of link capacity loss on SCSA include (a) less packets could be serviced without polling due to the shrunk radio budget, which can be explained by the "gate" below $1ms$ during the link capacity loss; (b) the algorithm becomes more vulnerable to bursty arrivals, which is visualized by the minor "spikes" climbing towards $3ms$ during the link capacity loss. Nevertheless, compared with UCSA, we can confirm that the connection level service variation due to particular link degradation is minimized by SCSA. Till now, the expected performance of SCSA in the second simulation scenario has been verified.

TABLE II
OVERALL NETWORK THROUGHPUT

| | UCSA ($Mbps$) | SCSA ($Mbps$) | improvement (%) |
|---|---|---|---|
| scenario one | 52.445344 | 53.590544 | 2.1836 |
| scenario two | 50.756944 | 51.012584 | 0.5037 |

*3) Network throughput:* The overall network throughputs measured in scenario one and scenario two, by UCSA and SCSA, are recorded in Table II. It is shown that SCSA improves network throughput over UCSA by $2.18\%$ in scenario one and $0.50\%$ in scenario two. This throughput improvement

---

[3]In the simulation we have observed no service provided to nrtPS connection 69 and 70 at SS8 from 2.0 second to 4.0 second. This implies that rtPS connections 67 and 68 at the same SS are also serviced less promptly and thus the connection queues are gradually built-up, during this time period. However, since connections 67 and 68 dominant the outbound transmission during this short-term link capacity loss, buffer overflow of these two connections was not found in the simulation.

quantitatively interprets the signaling overhead reduction capability of SCSA stated in Section IV. Considering the much larger number of connections running in practical systems than in the simulation experiments, SCSA may create appreciable revenue points for large-scale commercial deployments.

## VII. CONCLUSION

In this study we first presented a brief overview of the new IEEE 802.16 specifications on broadband wireless access networks. Then we proposed a new QoS control protocol design for single-carrier PMP mode wirelessMAN applications, and evaluated its performance by simulations. This proposed SCSA scheme enables each connection's contracted QoS parameters to control the service provided to the connection, which ensures the per-connection QoS guarantee. With some functionalities relocated from the base station to each subscriber station, signaling overhead is reduced. Moreover, by MAC-PHY cross-layer design for resource allocation, the new QoS control scheme is robust against wireless link degradation at particular subscriber station. Specifically, this cross-layer design averages the risk of particular channel condition deterioration over the entire network, and thereby the connection level service variation is minimized. Simulation results have firmly verified the expected performance of the new scheme.

## REFERENCES

[1] S. J. Vaughan-Nichols, "Achieving wireless broadband with WiMax," *Computer*, vol. 37, no. 6, pp. 10–13, June 2004.

[2] G. Goth, "Wireless MAN standard signals next-gen opportunities," *IEEE Distributed Systems Online*, vol. 5, no. 8, pp. 4–4, Aug. 2004.

[3] *IEEE Standard for Local and Metropolitan Area Networks Part 16: Air Interface for Fixed Broadband Wireless Access Systems*, IEEE Std 802.16-2004, Oct. 2004. [Online]. Available: http://ieeexplore.ieee.org/servlet/opac?punumber=9349

[4] *IEEE Std. 802.16-2001 IEEE Standard for Local and Metropolitan area networks Part 16: Air Interface for Fixed Broadband Wireless Access Systems*, IEEE Std 802.16-2001, Apr. 2002. [Online]. Available: http://ieeexplore.ieee.org/servlet/opac?punumber=7832

[5] *IEEE Standard for Local and metropolitan area networks — Part 16: Air Interface for Fixed Broadband Wireless Access Systems— Amendment 2: Medium Access Control Modifications and Additional Physical Layer Specifications for 2-11 GHz*, IEEE Std 802.16a-2003, Jan. 2003. [Online]. Available: http://ieeexplore.ieee.org/servlet/opac?punumber=8508

[6] K. Wongthavarawat and A. Ganz, "Packet scheduling for QoS support in IEEE 802.16 broadband wireless access systems," *Int. J. Commun. Syst.*, vol. 16, pp. 81–96, 2003.

[7] D. Niyato and E. Hossain, "Queue-aware uplink bandwidth allocation and rate control for polling service in IEEE 802.16 broadband wireless networks," vol. 5, no. 6, pp. 668–679, June 2006.

[8] H. S. Alavi, M. Mojdeh, and N. Yazdani, "A quality of service architecture for IEEE 802.16 standards," in *IEEE Asia-Pacific Conference on Communications*, Perth, Western Australia, Oct. 2005, pp. 249–253.

[9] J. Chen, W. Jiao, and H. Wang, "A service flow management strategy for IEEE 802.16 broadband wireless access systems in TDD mode," in *Proc. IEEE International Conference on Communications (ICC'05)*, vol. 5, Seoul, Korea, May 2005, pp. 3422–3426.

[10] S. A. Xergias, N. Passas, and L. Merakos, "Flexible resource allocation in IEEE 802.16 wireless metropolitan area networks," in *The 14th IEEE Workshop on Local and Metropolitan Area Networks (LANMAN'05)*, Chania, Greece, Sept. 2005, pp. 1–6.

[11] H. Wang, W. Li, and D. P. Agrawal, "Dynamic admission control and qos for 802.16 wireless MAN," in *IEEE Wireless Telecommunications Symposium (WTS'05)*, Pomona, CA, Apr. 2005, pp. 60–66.

[12] Q. Liu, S. Zhou, and G. B. Giannakis, "Queuing with adaptive modulation and coding over wireless links: Cross-layer analysis and design," *IEEE Trans. Wireless Commun.*, vol. 4, no. 3, pp. 1142–1153, May 2005.

[13] A. Scaglione and M. van der Schaar, "Cross-layer resource allocation for delay-constrained wireless video transmission," in *Proc. IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP'05)*, vol. 5, Philadelphia, PA, Mar. 2005, pp. v/909–v/912.

[14] S. Ramachandran, C. W. Bostian, and S. F. Midkiff, "A link adaptation algorithm for IEEE 802.16," in *Proc. IEEE Wireless Communications and Networking Conference (WCNC'05)*, vol. 3, New Orleans, LA, Mar. 2005, pp. 1466–1471.

[15] Y. Yao and J. Sun, "Study of UGS grant synchronization for 802.16," in *Proc. The Ninth International Symposium on Consumer Electronics (ISCE'05)*, Macau, China, June 2005, pp. 105–110.

[16] L. Zhang, "VirtualClock: A new traffic control algorithm for packet-switched networks," *ACM Transactions on Computer Systems*, vol. 9, no. 2, pp. 101–124, May 1991.

[17] The network simulator-ns-2. [Online]. Available: http://www.isi.edu/nsnam/ns/

[18] X. Bai, A. Shami, K. A. Meerja, and C. Assi, "New distributed QoS control scheme for IEEE 802.16 wireless access networks," in *Proc. IEEE Global Telecommunications Conference (GLOBECOM'06)*, San Francisco, CA, Nov.