

Scheduling Problems and Solutions in WiMAX Networks

Jia-Ming Liang, You-Chiun Wang, and Yu-Chee Tseng

Abstract—WiMAX is developed to support large-scale wireless broadband access. Defined in a series of IEEE 802.16 standards, three architectures of WiMAX networks are proposed to adapt to various environments. The point-to-multipoint architecture is used to manage a few number of devices, where each device is able to directly communicate with the central base station. The relay architecture deploys some special stations to act as intermediary between devices and the base station, where devices can choose whether or not to communicate with the base station through these intermediary. The mesh architecture is deployed to cover a large area, where all devices/stations are organized in an ad hoc fashion. Different challenges are arisen under different network architectures, leading the WiMAX scheduling problem to attract lots of research focus. In this chapter, we discuss the scheduling problems and their solutions under the three architectures of WiMAX networks, which covers the issues of how to improve network throughput, how to guarantee quality of service, and how to reduce energy consumption of devices. The comparison of these scheduling solutions is also given in the chapter.

Index Terms—IEEE 802.16, mesh, OFDM/OFDMA, point-to-multipoint, relay, resource management, WiMAX.



1 INTRODUCTION

WiMAX is an emerging wide-range wireless access technology for solving the last-mile communication problem, bridging the Internet and wireless local-area networks, and supporting broadband multimedia communication services [1], [2]. Recently, WiMAX networks have been widely deployed in many countries such as South Korea, India, and South Africa to provide low-cost Internet access [3], [4]. A series of IEEE 802.16 standards are defined to regulate WiMAX to support high-speed Internet access over long distances. Two types of accessing techniques, namely *orthogonal frequency division multiplexing (OFDM)* and *orthogonal frequency division multiple access (OFDMA)*, are employed in the WiMAX physical layer to realize the convergence of fixed and mobile broadband access through air interfaces. In a WiMAX network, the central *base station (BS)* is responsible for distributing the radio resource among *mobile subscriber stations (MSSs)* and scheduling the communication time of each MSS. To manage the resource, the standards define a *scheduler* in the *media access control (MAC)* layer of the BS but leave its detailed implementation as an open issue to provide the flexibility for the hardware manufacturers and network operators.

Depending on the application requirements and the covered areas, WiMAX defines three types of network architectures: 1) The *point-to-multipoint (PMP)* architecture consists of one BS and multiple MSSs, where each MSS can directly communicate with the BS. Such an architecture could be applied in those areas with sparse MSSs such as suburbs. 2) Under the *relay* architecture, several *relay stations (RSs)* are deployed to help relay the data between the BS and MSSs. Each MSS can choose either one-hop or two-hop (via an RS) communication to reach the BS. The relay architecture could be adequate to those areas with dense MSSs such as downtowns. 3) Under the *mesh* architecture, all *subscribe*

stations (SSs) are organized in an ad hoc fashion and each SS can reach the BS through a multihop manner. Compared to the above two architectures, the mesh architecture is usually adopted to cover a huge area such as metropolis or large islands. Explicitly, different architectures possess different network characteristics and constraints, which makes the WiMAX scheduling problem more challenging and interesting.

This chapter provides a comprehensive survey of the scheduling problems and solutions in WiMAX networks under the PMP, relay, and mesh architectures, which covers the following research issues:

- **Network throughput:** Since the objective of WiMAX is to provide broadband network access, we will introduce several scheduling schemes that target at improving network throughput. The concepts of control overhead reduction and concurrent transmissions will be adopted to help enhance throughput.
- **Quality of service (QoS):** The IEEE 802.16 standards classify all traffics into five QoS categories, each possessing different bandwidth requirements and delay constraints. We will discuss how to schedule MSSs' traffics so that their demands can be satisfied.
- **Energy consumption:** Communication is an energy-costly operation for mobile devices. We will survey some research efforts that adaptively adjust the communication powers of MSSs to balance between their energy consumption and the overall network throughput.

The rest of this chapter is organized as follows: Some background knowledge of WiMAX networks are given in the next section. Sections 3, 4, and 5 present the scheduling solutions for WiMAX networks under the PMP, relay, and mesh architectures, respectively. Section 6 concludes this chapter.

2 WiMAX NETWORKS

Below, we give an overview of WiMAX networks, which covers the topics of network architectures, accessing techniques in the physical layer, frame structures, and QoS service classes.

J.-M. Liang and Y.-C. Tseng are with the Department of Computer Science, National Chiao-Tung University, Hsin-Chu, 30010, Taiwan.

E-mail: {jmliang, wangyc, yctsen}@cs.nctu.edu.tw

Y.-C. Wang is with the Department of Computer Science and Engineering, National Sun Yat-sen University, Kaohsiung, 80424, Taiwan.

E-mail: ycwang@cse.nsysu.edu.tw

2.1 Network Architecture

To make the deployed networks be able to meet the application requirements or constraints imposed by the covered areas, WiMAX supports three types of network architectures, which are specified in different versions of IEEE 802.16 standards.

PMP architecture: Specified in the IEEE 802.16d and 802.16e standards [5], [6], PMP is a fundamental network architecture to support the wireless backhaul that enables high-speed Internet access (up to 70 Mbps) over long distances (up to 30 miles). Under the PMP architecture, the central BS can directly communicate with MSSs within its signal coverage, as shown in Fig. 1(a). In this case, the network will form a star topology centered at the BS. Those MSSs near the BS can receive stronger signals so that they could enjoy higher communication rates. On the other hand, those MSSs near the coverage boundary (such as MSS₁ and MSS₄) may receive weak signal power from the BS. Thus, they are asked to transmitted/received using lower communication rates so that more radio resource will be wasted. In addition, interfered by obstacles such as high buildings, trees, and mountains, the communication signal between the BS and an MSS would be weakened or even obstructed. This is called a *shadowing effect*. In this case, there could exist some *coverage holes* inside the BS's signal coverage and MSSs could not be able to communicate with the BS when they move into these coverage holes. Fig. 1(a) gives an example, where there is a coverage hole caused by the shadowing effect from the tree. MSS₅ may not receive the signal from the BS when it moves into the coverage hole.

Relay architecture: To improve network performance and solve the shadowing problem under the PMP architecture, the IEEE 802.16j standard [7] suggests deploying some RSs to help relay data between the BS and MSSs, as shown in Fig. 1(b). Each RS can be viewed as an 'extended' BS to enhance the received signal power at MSSs (such as MSS₁ and MSS₄) and eliminate the shadowing effect (such as MSS₅). The standard defines two types of RSs. When MSSs are not aware of the existence of RSs, these RSs are called *transparent*. Otherwise, they are *non-transparent*. Transparent RSs are used to increase network performance while non-transparent RSs are used to expand the BS's signal coverage. Transparent RSs are not responsible for arranging the radio resource to MSSs (such a job is handled by the BS), so they are easier to implement than the non-transparent RSs. Thus, this chapter aims at relay networks with transparent RSs. In a relay network, each MSS can choose to directly communicate with the BS or ask an RS to relay its data in a two-hop manner. However, any two MSSs or any two RSs cannot directly communicate with each other. In this way, the network will form a two-level tree rooted at the BS. Note that with RSs, concurrent RS-MSS communications may be realized due to spatial reuse.

Mesh architecture: Unlike the above two architectures, a mesh network consists of one BS and multiple *static* SSs (for example, these SSs can be set on the top of buildings to provide wireless access of the whole buildings). Specified by the IEEE 802.16d standard, all SSs will be organized in an ad hoc manner to cover a huge area. Two SSs can communicate with each other if they are within each other's transmission range. Each SS can act as either an

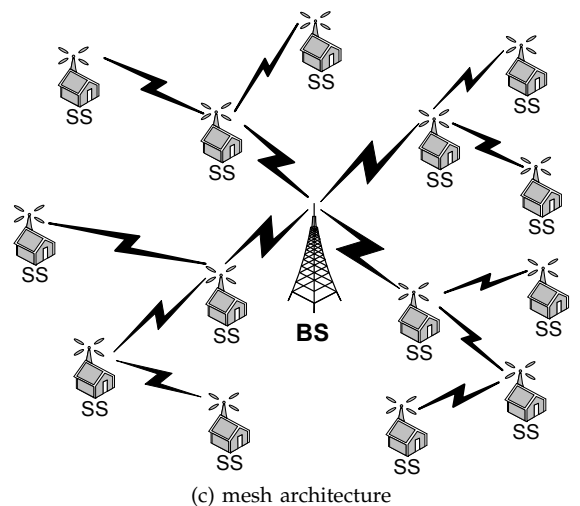
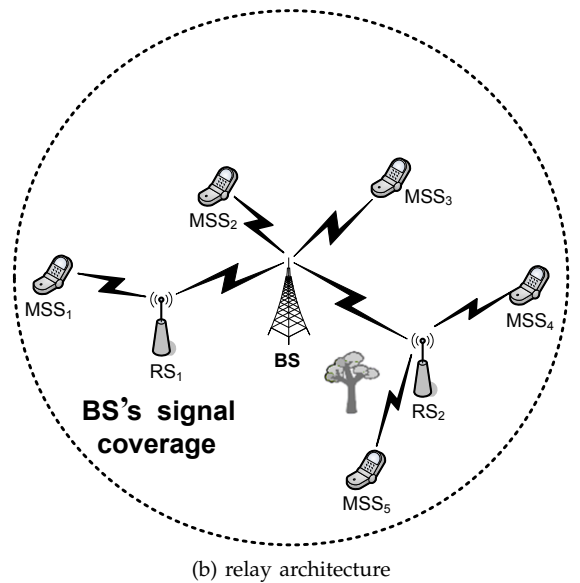
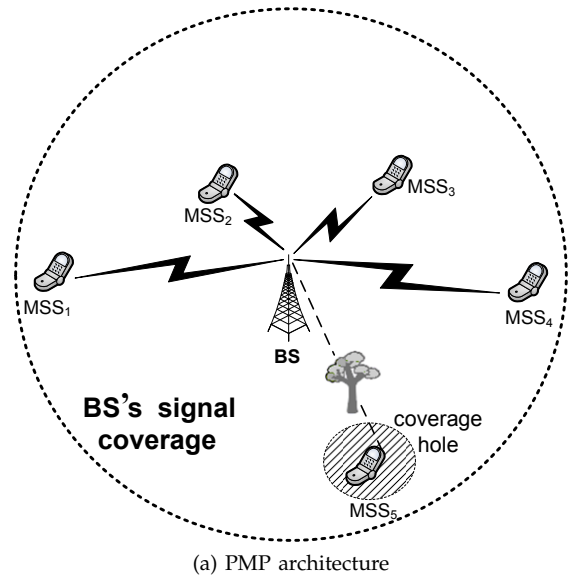


Fig. 1: The three network architectures supported by WiMAX: (a) Under the PMP architecture, the network will form a star topology and there could exist some coverage holes inside the BS's signal coverage. (b) Under the relay architecture, the network will form a two-level tree for communication purpose, where RSs help relay data between the BS and MSSs. (c) Under the mesh architecture, the BS constructs a routing tree for SSs to transmit/receive their data.

end point or a router to relay data for its neighbors. Since the BS is responsible for managing the radio resource, all

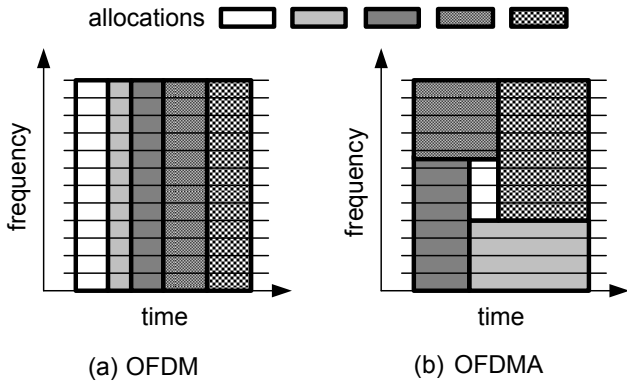


Fig. 2: Two accessing techniques adopted in the WiMAX physical layer, where the radio resource is distributed among five allocations. (a) Using OFDM, each SS has the full control of all subcarriers at different times. (b) Using OFDMA, different MSSs are allowed to access different subcarriers at the same time.

SSs have to send their requests containing traffic demands to the BS. Then, the BS will use the topology information along with SSs' requests to construct a routing tree for SSs to transmit/receive their data, as shown in Fig. 1(c). It can be observed that more concurrent communications could coexist since some SSs are deployed far away from each other.

2.2 Accessing Techniques in The Physical Layer

The WiMAX physical layer supports two types of accessing techniques, OFDM and OFDMA, as shown in Fig. 2.

OFDM technique: The mesh architecture adopts OFDM as the accessing technique in the physical layer. OFDM supports *non-line of sight (NLOS)* communications and multicarrier transmissions, where each SS is given the complete control of all subcarriers. The BS adopts the concept of *time division multiple access (TDMA)* to share the radio resource among all SSs. In other words, for multiple SSs that are within each other's transmission range, only one SS is allowed to access the channel at any time. Therefore, the BS only needs to determine which time slot should be allocated to which SS. Fig. 2(a) gives an example, where the radio resource is distributed among into five allocations. Each allocation can be viewed as a rectangle whose height covers all available frequency bands. Any two allocations do not overlap in the time domain.

OFDMA technique: The PMP and relay architectures adopt OFDMA as the accessing technique in the physical layer to support the mobility of MSSs. Unlike OFDM, different MSSs are allowed to transmit/receive data through different subcarriers at the same time to enhance the signal power of the MSS. Fig. 2(b) gives an example, where the five allocations together constitute the whole radio resource. Since the BS needs to determine which time slot and which subcarrier should be allocated to which MSS, an OFDMA BS will be more complex than an OFDM BS.

Note that a scheduler only determines the sizes of allocations but does not take care of how to arrange these allocations to fit into the two-dimensional time-frequency array (in Fig. 2). Such an issue has been addressed in the studies of [8]–[10].

2.3 Frame Structures

In WiMAX networks, the radio resource is divided into *frames*. According to different network architectures, vari-

TABLE 1: The six MCSs supported by WiMAX: Using different MCS levels, each slot can carry different amount of data and each MCS requires a minimum *signal to interference plus noise ratio (SINR)*. A higher MCS level requires a higher SINR and can carry more data. On the contrary, a lower MCS level requires a lower SINR and can carry less data.

level	MCS	data carried by each slot	minimum SINR
1	QPSK 1/2	48 bits	6 dBm
2	QPSK 3/4	72 bits	8.5 dBm
3	16QAM 1/2	96 bits	11.5 dBm
4	16QAM 3/4	144 bits	15 dBm
5	64QAM 2/3	192 bits	19 dBm
6	64QAM 3/4	216 bits	21 dBm

ous frame structures are also defined:

PMP architecture: Since the PMP architecture adopts the OFDMA accessing technique, the frame will be a two-dimensional array with time units in the time domain and subchannels in the frequency domain, as shown in Fig. 3(a). The basic unit of a frame is called a *subchannel-time slot* (or simply *slot*). Each frame is further divided into a *downlink subframe* and an *uplink subframe*. A downlink subframe is composed of the *preamble*, *control*, and *data* portions, while an uplink subframe only has the data portion. The preamble portion is used for time synchronization. The control portion contains the *frame control header (FCH)*, *downlink map (DL_MAP)*, and *uplink map (UL_MAP)* fields. The DL_MAP and UL_MAP fields are used to indicate the downlink and uplink resource allocation in the current frame, respectively. In the data portion, each allocation is a subarray of slots, called a *burst*. From Fig. 3(a), each burst in the downlink subframe is shaped by a rectangle whose width may be multiple subchannels. On the other hand, the bursts in the uplink subframe should be arranged in a row-wise manner, where each burst has a width of only one subchannel. In practice, each MSS can be allocated with more than one burst. However, any two bursts cannot overlap with each other.

Each downlink/uplink burst is with a *modulation and coding scheme (MCS)* and requires one *information element (IE)* recorded in the DL_MAP/UL_MAP field to indicate its size and location in the downlink/uplink subframe. Table 1 lists the six MCSs support by WiMAX. Note that each burst can only carry the data of exact one MSS. Therefore, the number of bursts (and thus IEs) will increase when the BS admits more MSSs to access the radio resource. Each IE requires 60 bits encoded by QPSK1/2 (that is, the lowest MCS level). From Table 1, each slot can carry data of 48 bits, so an IE will occupy $\frac{5}{4}$ slots. Because IEs and bursts share the same space in the downlink subframe, too many IEs may degrade the network performance.

Relay architecture: Since both the PMP and relay architectures adopt the OFDMA accessing technique, their frame structures will share some common features. For example, the frame is also modeled by a two-dimensional array over both the time and frequency domains. The bursts allocated in the downlink subframe are shaped by rectangles with different widths while the bursts in the uplink subframe are arranged in a row-wise manner. In addition, each downlink/uplink burst spends one IE in the DL_MAP/UL_MAP field to record its corresponding allocation information.

However, because of the existence of RSs, there are two types of frames, namely *BS frames* and *RS frames*. Generally speaking, a BS frame has a 'complementary' RS frame, as

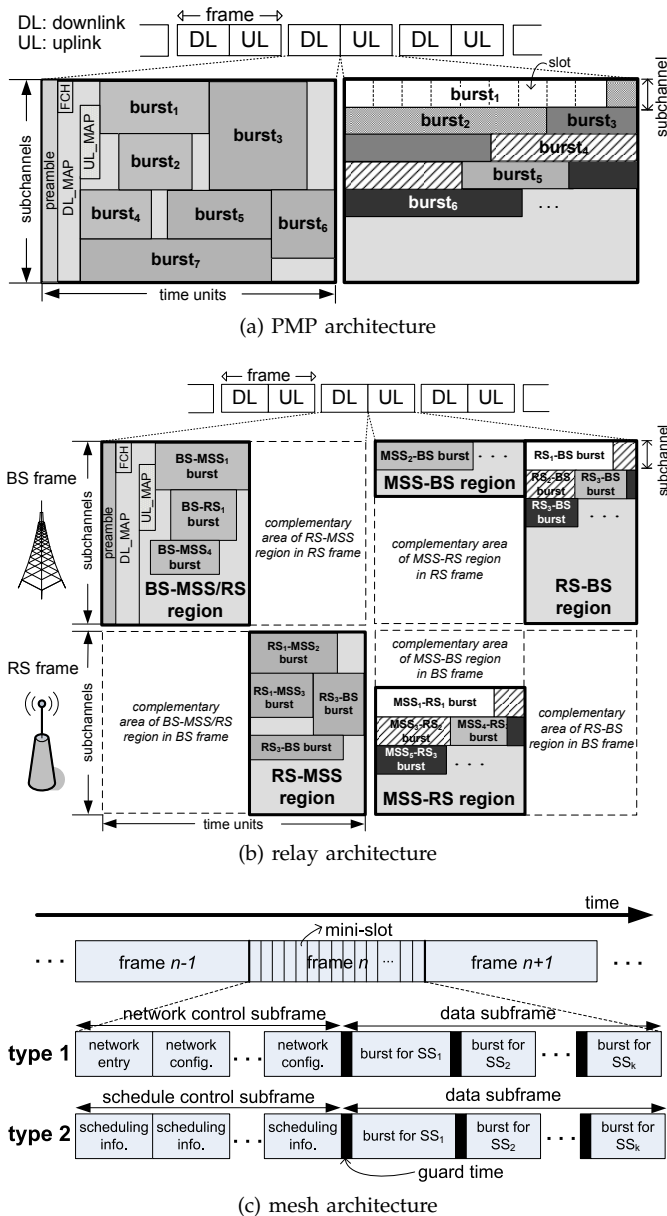


Fig. 3: The frame structures under different network architectures: Under the PMP and relay architectures, the frame is modeled by a two-dimensional array over both the time and frequency domains. On the other hand, under the mesh architecture, the frame is modeled by an one-dimensional array over the time domain.

shown in Fig. 3(b). For a BS frame, its downlink subframe has a *BS-MSS/RS region* to allocate downlink bursts for the BS to transmit data to MSSs or RSs; its uplink subframe has an *MSS-BS region* and an *RS-BS region* to allocate uplink bursts for MSSs and RSs to submit their data to the BS, respectively. On the other hand, for an RS frame, its downlink subframe has an *RS-MSS region* to allocate downlink bursts for the RS to relay data to MSSs; its uplink subframe has an *MSS-RS region* to allocate uplink bursts for MSSs to submit their data to the BS through the RS. Each RS is considered as ‘bufferless’ in the sense that the data received by the RS from the BS/MSS must be delivered to the MS/BS during the same frame. Taking the uplink subframe in Fig. 3(b) as an example, since an MSS_1 - RS_1 burst is allocated in the *MSS-RS region*, there must be an RS_1 - BS burst allocated in the *RS-BS region*.

Because the BS is the only receiver, any two bursts in the *BS-MSS/RS*, *MSS-BS*, and *RS-BS* regions cannot overlap

with each other. However, by exploiting spatial reuse, concurrent *MSS-RS* or *RS-MSS* communications may be allowed. Therefore, some bursts could be overlapped with each other in the *RS-MSS* and *MSS-RS* regions to improve network efficiency.

Mesh architecture: Taking OFDM as the accessing technique in the physical layer, the frame under the mesh architecture is modeled by an one-dimensional array over the time domain. The basic unit of each frame is called a *mini-slot*. Two types of frames are defined, as shown in Fig. 3(c). A type-1 frame consists of a *network control subframe* and a *data subframe*, where the former carries some network formation information such as how to construct the routing tree while the latter carries the bursts of SSs. The length of the network control subframe is fixed. For each burst, it requires a *guard time* in front of it to conduct time synchronization and avoid propagation delay interfering the following transmission. Such a guard time is usually viewed as transmission overhead because it does not carry the SS’s data. Note that the burst of each SS may mix its downlink and uplink data. On the other hand, a type-2 frame has a fixed-length *schedule control subframe* used to specify the resource allocation in the following data subframe. Each *scheduling information* field contains the burst accessing information such as which mini-slots in the corresponding burst are used for uplink or downlink communication.

Type-1 frames are used for network configurations and type-2 frames are used for normal transmission. It can be observed that the transmission overhead caused by guard times will degrade network performance and thus how to alleviate these overhead is a critical issue.

2.4 QoS Service Classes

To satisfy the different requirements of various data traffics, WiMAX defines five types of QoS service classes:

Unsolicited grant service (UGS): The UGS class provides fixed periodic bandwidth allocation for *constant bit rate (CBR)* traffics such as E1/T1 circuit emulation. Each MSS or SS only needs to negotiate with the BS about the QoS parameters such as maximum sustained rate, maximum latency, and tolerated jitter at the first time when the connection is established. Then, no further negotiation is required. The UGS class can guarantee the maximum latency for those delay-critical real-time services. However, the radio resource may be wasted if the granted traffics do not fully utilize the allocated bandwidth.

Real-time polling service (rtPS): The rtPS class supports *variable bit rate (VBR)* traffics such as compressed videos. Unlike UGS, the BS has to periodically poll each MSS or SS for its QoS parameters such as maximum sustained rate, maximum latency, tolerated jitter, and minimum reserved rate. The benefit is that the BS can adjust bandwidth allocation according to the real demands of traffics. However, periodical polling may also spend the radio resource.

Extended real-time polling service (ertPS): The ertPS class is specially designed for *voice over IP (VoIP)* with silence suppression, where no traffic is sent during silent periods. Both ertPS and UGS share the same QoS parameters. The BS will allocate the bandwidth with the maximum sustained rate when the VoIP traffic is active and no bandwidth when it becomes silent. In this way, the

BS only has to poll MSSs or SSs during the silent period to determine whether their VoIP traffics become active again.

Non-real-time polling service (nrtPS): The nrtPS class considers those non-real-time traffics with minimum reserved rates. The *file transfer protocol (FTP)* is one representative example. The BS will preserve bandwidth according to the minimum reserved rate to avoid starving the non-real-time traffic.

Best effort service (BE): All other traffics belong to this service class. The BS will distribute the remaining bandwidth (after allocating to the traffics of all other four service classes) to the traffics of the BE class, so there is no guarantee of throughput or delay.

Table 2 summarizes the notations used in this chapter.

3 SCHEDULING SOLUTIONS UNDER THE PMP ARCHITECTURE

Under the PMP architecture, the BS is responsible for managing the radio resource for all MSSs around it. The major objective is to guarantee QoS requirements of MSSs and improve network throughput. Below, we introduce three scheduling solutions where the BS deals with scheduling in the basis of connections, MSSs, and subchannels.

3.1 Connection-based Scheduling Solution

The work of [11] schedules *connections* according to their traffic types, where each MSS may contain one or multiple connections. For UGS and ertPS connections, the BS always allocates a fixed amount of resource to them. Then, the remaining resource is allocated to other connections according to their priorities, which are calculated as follows:

- **rtPS connections:** For each connection j , we adopt an indicator $R_j(t)$ to evaluate its channel quality at time t , which can be calculated by the number of packets carried by a time slot under that channel quality. Let R_N denote the maximum value of $R_j(t)$, $\forall j, t$. Then, the priority of an rtPS connection i is defined by

$$\phi_i(t) = \begin{cases} \beta_{rt} & \text{if } R_i(t) > 0 \text{ and } F_i(t) < 1 \\ \beta_{rt} \times \frac{R_i(t)}{R_N} \times \frac{1}{F_i(t)} & \text{if } R_i(t) > 0 \text{ and } F_i(t) \geq 1 \\ 0 & \text{if } R_i(t) \leq 0, \end{cases} \quad (1)$$

where $\beta_{rt} \in [0, 1]$ is a coefficient to evaluate the priority of rtPS connections and $F_i(t)$ is an indicator to measure the delay of connection i :

$$F_i(t) = T_i - \Delta T_i - W_i(t) + 1,$$

where $\Delta T_i \in [0, T_i]$ is the guard-time region ahead of connection i 's deadline T_i and $W_i(t) \in [0, T_i]$ is the longest packet waiting time of connection i . In Eq. (1), when $F_i(t) < 1$ under a positive $R_i(t)$, the highest priority β_{rt} in the rtPS class is given to connection i since its packet deadline is approaching (that is, $W_i(t) \in (T_i - \Delta T_i, T_i]$). Otherwise, the priority of connection i will be proportional to its channel quality $R_i(t)$ and inverse proportional to the delay indicator $F_i(t)$. Explicitly, when $R_i(t)$ is zero, which means that the channel quality of connection i is too bad to transmit data, a zero priority is set to let the BS neglect connection i in the current frame.

- **nrtPS connections:** Similar to the rtPS class, the priority of each nrtPS connection i is defined by

$$\phi_i(t) = \begin{cases} \beta_{nrt} & \text{if } \hat{\eta}_i(t) < \eta_i \text{ and } R_i(t) > 0 \\ \beta_{nrt} \times \frac{R_i(t)}{R_N} \times \frac{\eta_i}{\hat{\eta}_i(t)} & \text{if } \hat{\eta}_i(t) \geq \eta_i \text{ and } R_i(t) > 0 \\ 0 & \text{if } R_i(t) \leq 0, \end{cases} \quad (2)$$

where $\beta_{nrt} \in [0, 1]$ is a coefficient to evaluate the priority of nrtPS connections and $\hat{\eta}_i(t)$ and η_i are the average transmission rate and minimum reserved rate of connection i , respectively. In Eq. (2), when $\hat{\eta}_i(t) < \eta_i$ under a positive $R_i(t)$, which means that connection i is at the risk of being starved, the highest priority β_{nrt} in the nrtPS class should be given to connection i . Otherwise, when more resource is received by connection i (that is, larger $\hat{\eta}_i(t)$), a lower priority is set to maintain a certain degree of fairness.

- **BE connections:** For each BE connection, its priority is defined by

$$\phi_i(t) = \beta_{BE} \times \frac{R_i(t)}{R_N},$$

where $\beta_{BE} \in [0, 1]$ is a coefficient to evaluate the priority of BE connections. It can be observed that a connection with better channel condition will be given a higher priority.

Since the rtPS class possesses a strict delay constraint and the BE class has no QoS concern, it is suggested to set $\beta_{rt} > \beta_{nrt} > \beta_{BE}$.

Sorting by their priorities in a decreasing order, the BS selects connections to serve in sequence. A simulation with two rtPS connections and four nrtPS connections is conducted, and the results show that the delay outage probability¹ of the rtPS connections is below 5% and the average transmission rate of all connections can reach 6 Mbps. Thus, using priorities, those connections with the strict delay constraint or better channel condition will be served first to guarantee their QoS requirements and improve network throughput. However, since the BS conducts scheduling in a connection-based manner, it would generate many IEs in DL_MAP or UL_MAP and thus may hurt system performance. This issue will be addressed in the next section.

3.2 MSS-based Scheduling Solution

To reduce both scheduling complexity and IE overhead, the work of [12] suggests scheduling MSSs rather than connections. In addition, the data of each MSS is divided into *real-time* and *non-real-time* ones, where real-time data contains those data in UGS, rtPS, and ertPS service classes while non-real-time data contains those data in nrtPS and BE service classes. The idea is to limit the amount of real-time data to be scheduled so that network throughput can be improved by giving more resource to the non-real-time data of those MSSs with good channel quality.

Given the channel rate c_i , the amount of buffered real-time data b_i^R , the real-time data rate r_i^R , and the non-real-time satisfaction ratio s_i^N of each MSS i , the BS will assign

1. The delay outage probability is the probability that packets miss their deadlines.

TABLE 2: Summary of notations.

notation	definition
$\phi_i(t)$	the priority of connection i at time t
$\beta_{rt}/\beta_{nrt}/\beta_{BE}$	the coefficients to evaluate the priority of rtPS/nrtPS/BE connections
$R_i(t)$	the quality of the i th channel at time t
R_N	the maximum value of $R_i(t)$
$F_i(t)$	the indicator to measure the delay of connection i
ΔT_i	the guard-time region ahead of connection i 's deadline T_i
$W_i(t)$	the longest packet waiting time of connection i at time t
$\hat{\eta}_i(t)$	the average transmission rate of connection i at time t
η_i	the minimum reserved rate of connection i
n	the number of total MSSs
p_i	the priority of MSS i
c_i	the current channel rate of MSS i
c_i^{avg}	the average channel rate of MSS i
b_i^R/b_i^N	the amount of buffered real-time/non-real-time data of MSS i
r_i^R/r_i^N	the real-time/non-real-time data rates of MSS i
s_i^N	the non-real-time satisfaction ratio of MSS i
f_T	the window size used for observation (in frames)
$a_{i,j}^N$	the amount of resource allocated to MSS i 's non-real-time data at frame j
d_i	the amount of urgent data of MSS i
a_i	the amount of data allocated by the BS to MSS i
σ_i	the summation of allocated slots for all MSSs except MSS i
F	the number of free slots in the current frame
δ	a ratio to serve MSSs' real-time data
P_i^k	the priority for MSS i on subchannel k
S_i^k	the channel state of MSS i on subchannel k
Q_i	the QoS satisfaction indicator of MSS i
Q_{ij}	the QoS satisfaction indicator of MSS i 's connection j
K	the total number of subchannels
b_i^k	the number of bits that can be carried by one subcarrier of MSS i in one OFDMA symbol on subchannel k
b_i^{\max}	the maximum bits per symbol carried by one subcarrier
b_{\max}^k	the maximum value of b_i^k ($i = 1..n$)
d_k	the normalized deviation of channel quality for channel k
w_{ij}, f_{ij}	the longest and the maximum tolerable waiting times of packets in connection j of MSS i , respectively
q_{ij}, μ_{ij}	the queue length and its threshold of connection j of MSS i , respectively
q_i	the queue length of SS i
h_i	the hop count from the BS to SS i
h_{\max}	the maximum hop count in the network
\hat{S}	a set of SSs selected for concurrent transmissions
H	the minimum hop count to allow SSs to concurrently transmit

it with a priority

$$p_i = c_i \times \frac{c_i}{c_i^{avg}} \times \frac{b_i^R}{r_i^R} \times \frac{1}{s_i^N}, \quad (3)$$

where c_i^{avg} is the average channel rate of MSS i and

$$s_i^N \triangleq \min \left\{ 1, \frac{\sum_{j=0}^{f_T-1} a_{i,f_c-j}^N}{f_T \times r_i^N} \right\},$$

where f_T is the window size used for observation (in frames) and a_{i,f_c-j}^N is the amount of resource allocated to MSS i 's non-real-time data at frame $(f_c - j)$. In Eq. (3), the first term c_i means that an MSS with better channel quality will be given a higher priority to improve network throughput. The second term c_i/c_i^{avg} is to raise the priorities of those MSSs that encounter chronic bad channel conditions to avoid starving them. The third term b_i^R/r_i^R and the last term $1/s_i^N$ are used to give a higher priority for those MSSs that queue a lot of real-time data and possess lower non-real-time satisfaction ratios, respectively.

The BS then sorts MSSs by their priorities in a decreasing order and schedules them using this order:

- 1) For each MSS i with the amount of d_i urgent data, which is real-time data whose packets will be dropped if the BS does not schedule it in the current frame, the BS allocates it with the amount of resource $a_i = \min\{c_i \times (F - \sigma_i), d_i\}$, where F is the number of

free slots in the current frame and

$$\sigma_i = \sum_{\forall j, j \neq i} \left\lceil \frac{a_j}{c_j} \right\rceil,$$

is the summation of allocated slots for all MSSs except MSS i .

- 2) The BS then select the first $\lceil \delta n \rceil$ MSSs to serve their real-time data, where $0 < \delta < 1$ and n is the number of total MSSs. For each such MSS i , the BS allocates it with the amount of resource $a_i = \min\{c_i \times (F - \sigma_i), b_i^R\}$.
- 3) For each MSS i , the BS allocates it with the amount of resource $a_i = \min\{c_i \times (F - \sigma_i), b_i^R + b_i^N\}$.

It can be observed that in step 2 the BS does not serve the real-time data of *all* MSSs to reduce the IE overhead and prevent those real-time data with bad channel condition from occupying network resource. Therefore, non-real-time data with good channel condition can have an opportunity to transmit their packets to improve network throughput and maintain fairness. In addition, the BS conducts scheduling in an MSS-based manner so that the amount of IE overhead can be reduced, as compared with the connection-based scheduling scheme. A simulation with up to 70 MSSs is conducted, and the results show that the proposed scheme incurs no real-time packet dropping while guarantees non-real-time rate satisfaction.

3.3 Subchannel-based Scheduling Solution

The above two studies assume that all subchannels have the similar quality. The work of [13] considers that subchannels may have different qualities and thus schedules MSSs' connections in a subchannel-based manner. In particular, for each MSS i on subchannel k , the BS calculates a priority $P_i^k = S_i^k \times Q_i$, where S_i^k reflects the channel state of MSS i on subchannel k and Q_i is MSS i 's QoS satisfaction indicator. Here, the channel state is defined by

$$S_i^k = \frac{b_i^k}{b_{\max}} \times \frac{\sum_{k=1}^K (b_i^{\max} - b_i^k + 1)}{K b_i^{\max}}, \quad (4)$$

where b_i^k is the number of bits that can be carried by one subcarrier of MSS i in one OFDMA symbol on subchannel k , b_i^{\max} is the maximum value of b_i^k ($k = 1..K$), b_{\max} is the maximum bits per symbol carried by one subcarrier, and K is the total number of subchannels. In Eq. (4), the first term b_i^k/b_{\max} quantifies the normalized quality of the subchannel. The remaining part of Eq. (4) indicates the normalized deviation of channel quality of MSS i on different subchannels. A larger deviation means that some subchannels have lower quality and others have higher quality. In this case, if the BS gives a higher priority to MSS i , it may send data through its good subchannels. On the other hand, the QoS satisfaction indicator Q_i is the maximum connection priority of MSS i , where each connection j 's priority Q_{ij} is defined according to the service classes:

- **UGS and ertPS:** The largest value of Q_{ij} (for example, 2) is set to make UGS and ertPS connections have the highest priority to meet their delay requirements.
- **rtPS:** Let w_{ij} and f_{ij} denote the longest waiting time and the maximum tolerable waiting time of packets in connection j of MSS i , respectively. Then,

$$Q_{ij} = \begin{cases} \beta_{rt} & \text{if } w_{ij} \geq f_{ij} \\ \beta_{rt} \times \frac{w_{ij}}{f_{ij}} & \text{otherwise,} \end{cases}$$

where $\beta_{rt} \in [0, 1]$. It can be observed that when packet deadline is approaching, the highest connection priority in rtPS class is given; otherwise, the connection priority is proportional to the normalized queuing delay w_{ij}/f_{ij} .

- **nrtPS:** Let q_{ij} and μ_{ij} denote the current queue length and the length threshold of connection j of MSS i , respectively. Then,

$$Q_{ij} = \begin{cases} \beta_{nrt} & \text{if } q_{ij} \geq \mu_{ij} \\ \beta_{nrt} \times \frac{q_{ij}}{\mu_{ij}} & \text{otherwise,} \end{cases}$$

where $\beta_{nrt} \in [0, 1]$. It can be observed that when the connection is starved, the highest connection priority in nrtPS class is given; otherwise, the connection priority is proportional to the normalized queuing length q_{ij}/μ_{ij} .

- **BE:** The smallest values of Q_{ij} (for example, -1) is given since BE connections have no QoS concern.

On the other hand, for each subchannel k , the BS calculates its normalized deviation of channel quality as follows:

$$d_k = \frac{\sum_{i=1}^n (b_{\max}^k - b_i^k)}{n b_{\max}^k},$$

where b_{\max}^k is the maximum value of b_i^k ($i = 1..n$) and n is the total number of MSSs. Note that the above equation is

similar to the second term in Eq. (4). Then, the BS sorts all subchannels according to their d_k values in a decreasing order and selects these subchannels in sequence. For each selected subchannel k , the BS determines which MSS i can send data through the subchannel according to its priority P_i^k . An MSS with a higher priority can be served first. For each served MSS i , the BS sorts its connections according to their connection priorities Q_{ij} in a decreasing order and allocates resource to each connection to satisfy its minimum reserved rate in sequence.

After allocating resource to satisfy the minimum reserved rate of each connection, the BS will distribute the remaining resource among those MSSs with the best channel quality (on certain subchannels). By considering subchannel diversity, the BS can arrange MSSs to transmit/receive data through their suitable subchannels to improve network throughput. For the simulation study, [13] shows that in a small network (with 10 MSSs), the network throughput can approximate the maximum one and real-time traffics encounter lower packet dropping even when the channel quality becomes worse.

4 SCHEDULING SOLUTIONS UNDER THE RELAY ARCHITECTURE

Under the relay architecture, MSSs can select RSs to relay their data. With RSs, concurrent transmissions can be realized due to spatial reuse. Below, we introduce two scheduling solutions that use RSs to improve network throughput and reduce energy consumption of MSSs.

4.1 Scheduling Solution Using RSs to Improve Network Throughput

The studies of [14], [15] consider using spatial reuse to improve network throughput under the relay architecture. The idea is to first let each MSS select its 'best' path to reach the BS and then check whether some links can be transmitted simultaneously to improve network throughput. To help MSSs select their paths, a *cost* for each communication link is defined by the time of a bit transmitted by the highest MCS level of that link. In other words,

$$\text{cost} \triangleq 48 \text{ bits} \div (\text{bits carried by each slot using the highest MCS level of the link}).$$

For example, according to Table 1, the cost of a link is $48/72 = 2/3$ if the MCS level 2 (that is, QPSK 3/4) is used. Then, the cost of a path is defined by the sum of link costs along that path. Fig. 4 gives an example, where the costs of paths $p(MSS_1 - RS_1, RS_1 - BS)$, $p(MSS_1 - RS_2, RS_2 - BS)$ and $p(MSS_1, BS)$ are $(1/2 + 2/9) = 13/18$, $(1/4 + 1/3) = 7/12$, and 1, respectively. Then, for each MSS, it will select the path with the minimum cost as its path. For example, MSS 1 will select $p(MSS_1 - RS_2, RS_2 - BS)$ as its path.

After determining each MSS's path, those bursts for MSS-RS communications will be divided into multiple *transmission groups*, where all links in one transmission group are allowed to transmit simultaneously (by employing spatial reuse). In particular, sorting all MSS-RS links according to their costs in a decreasing order, we then select links following that order. For each selected link, we check whether or not it can join a transmission group (if this link will not cause interference to the existing links in

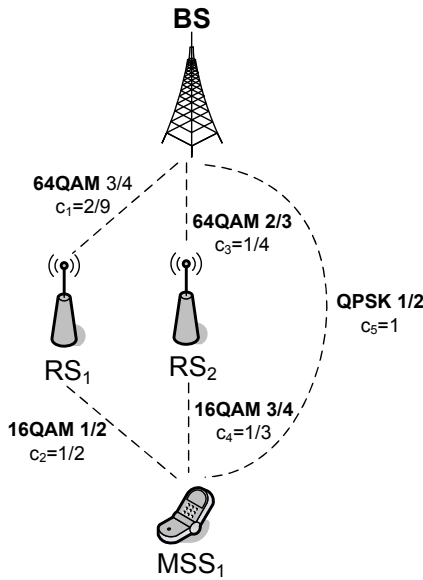


Fig. 4: An example of calculating path costs. The cost of a path is the sum of all link costs along that path, where the cost of each link is defined by the duration to transmit one bit using the highest MCS level of that link. Each MSS will select the path with the minimum cost as its path to the BS.

that transmission group). Here, the transmission groups are sorted by its maximum link cost in a decreasing order. If the link cannot join any transmission group, it will organize a transmission group containing itself. The above operation is repeated until all MSS-RS links are checked.

Allowing MSSs and RSs to use their highest MCS levels to transmit and adopting spatial reuse to make concurrent MSS-RS communications possible, the network throughput can be improved. A large-scale network with 300 MSSs is conducted in the simulation study. The results show that with RSs, the network throughput can be improved up to 85% compared to the case without RSs. However, both high MCS level and spatial reuse also make MSSs use large transmission powers for communications, which force them to consume more energy. This issue will be addressed in the next section.

4.2 Scheduling Solution Using RSs to Conserve MSSs' Energy

Given the uplink requests of MSSs, the work of [16] considers the *energy-conserved uplink resource allocation (EURA) problem* under the relay architecture, which determines how to allocate resource to MSSs such that their energy consumption is minimized under the constraint that MSSs' uplink requests are met (when resource is sufficient). EURA is proved to be NP-hard and a 2-phase heuristic is proposed to help the BS to arrange uplink bursts for MSSs in terms of their paths (either directly to the BS or via an RS), MCSs, and transmission powers:

Phase 1: The first phase tries to allocate the minimum amount of resource to satisfy MSSs' requests by employing spatial reuse to compactly arrange their bursts. To check whether spatial reuse can be employed, the BS calculates the *maximum tolerable interference (MTI)* for each MSS when it selects an RS to relay its data using one MCS level at its maximum power. Two uplink bursts can be transmitted simultaneously if the transmissions will not make the interference of the corresponding MSSs exceed

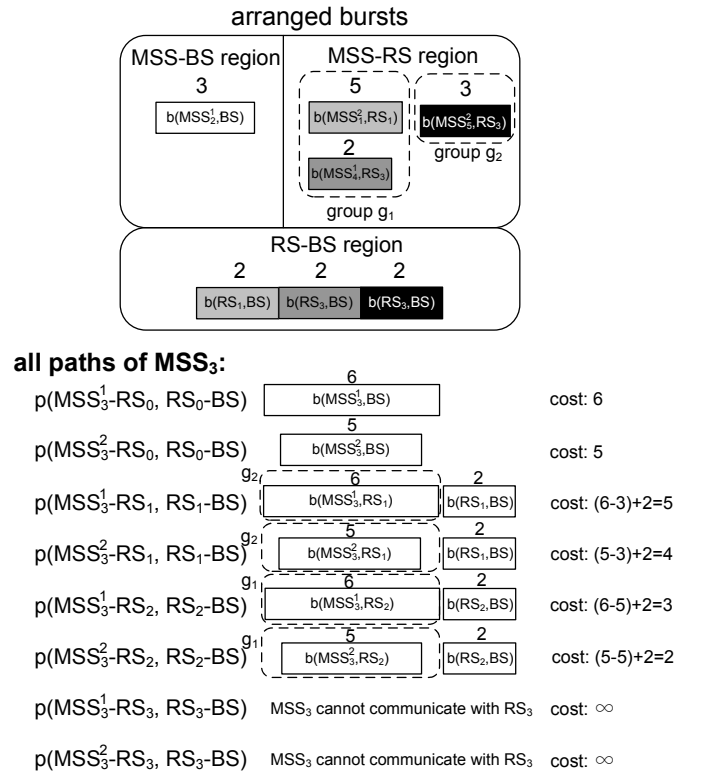


Fig. 5: An example of determining the path and corresponding burst(s) of each MSS in phase 1. Only the bursts in the MSS-RS region are allowed to employ spatial reuse to transmit simultaneously. The BS will calculate the costs of all possible paths for each MSS and select the one with the minimum cost as the MSS's path. In this example, the path $p(MSS_3^2 - RS_2, RS_2 - BS)$ is selected for MSS 3 and two corresponding bursts $b(MSS_3^2, RS_2)$ and $b(RS_2, BS)$ will be arranged.

their MTIs. In this case, these two bursts can be *grouped* together. Fig. 5 gives an example, where $b(MSS_i^k, BS)$, $b(MSS_i^k, RS_j)$, $b(RS_j, BS)$ denote the bursts for the communications between MSS i and the BS using MCS level k , between MSS i and RS j using MCS level k , and between RS j and the BS, respectively. It can be observed that only the bursts in the MSS-RS region can be grouped together.

Then, the BS evaluates the *cost* when arranging the uplink burst(s) for an MSS i using each MCS level. When MSS i directly communicates with the BS using MCS level k , the cost will be the length of burst $b(MSS_i^k, BS)$. On the other hand, when MSS i selects RS j for relay using MCS level k , the cost will be the sum of the length of burst $b(RS_j, BS)$ and the *extra length* arisen from burst $b(MSS_i^k, RS_j)$ if this burst joins a group. Fig. 5 gives an example. Consider the case when MSS 3 selects RS 1 as relay using MCS level 1 (in other words, MSS 3 selects the path $p(MSS_3^1 - RS_1, RS_1 - BS)$). Since the maximum burst length of group g_2 is only 3, if burst $b(MSS_3^1, RS_1)$ joins group g_2 , the extra length arisen from burst $b(MSS_3^1, RS_1)$ will be $6 - 3 = 3$. Consider another case when MSS 3 selects the path $p(MSS_3^2 - RS_2, RS_2 - BS)$. Since the maximum burst length of group g_1 (that is, 5) is equal to the length of burst $b(MSS_3^2, RS_2)$ (that is, 5), there is no extra length arisen from burst $b(MSS_3^2, RS_2)$. Note that when MSS i cannot communicate with RS j due to serious interference, the costs of those corresponding paths are set to infinity. After determining the costs of all possible paths for MSS i , the BS selects the path with the minimum cost (and arranges the corresponding bursts). If there is a tie, the BS will arbitrary select one path.

Phase 2: After determining the initial path, MCS, and transmission group of each MSS, the second phase checks whether it can further reduce the energy consumption of MSSs by lowering down their MCSs or changing the current paths to reduce their transmission powers. In particular, for each MSS i , suppose that it selects RS j as its relay using MCS level k and its corresponding burst $b(MSS_i^k, RS_j)$ belongs to group g_a in phase 1. Let $c(MSS_i, RS_j, k, g_a)$ and $e(MSS_i, RS_j, k, g_a)$ denote the cost and energy consumption of MSS i under the above situation, respectively. Then, the *energy-saving ratio* of MSS i is defined by

$$\frac{\Delta_E(MSS_i, RS_{j'}, k', g_a)}{\Delta_C(MSS_i, RS_{j'}, k', g_a)} = \frac{e(MSS_i, RS_j, k, g_a) - e(MSS_i, RS_{j'}, k', g_a)}{c(MSS_i, RS_{j'}, k', g_a) - c(MSS_i, RS_j, k, g_a)},$$

when MSS i changes to select RS j' as its relay using new MCS level k' and its corresponding burst $b(MSS_i^{k'}, RS_{j'})$ now joins to a new group g_a' . Among all possible combinations of (j', k', g_a') , the BS will select the one with the maximum ratio $\Delta_E(MSS_i, RS_{j'}, k', g_a')/\Delta_C(MSS_i, RS_{j'}, k', g_a')$ with $\Delta_E(MSS_i, RS_{j'}, k', g_a') \geq 0$ and then change the settings of MSS i accordingly. After changing the settings of all MSSs (if feasible), the BS will adjust the transmission power of each MSS according to its new path and MCS.

In sum, in phase 1 the BS tries to use the minimum amount of resource to satisfy the traffic demands of all MSSs and then in phase 2 it adjusts the paths and MCSs of all MSSs to reduce their energy consumption. In this way, not only the MSSs' uplink requests can be met but also their energy can be conserved. Simulation results show that when there are 8 RSs and 20 MSSs, the proposed scheme can save up to 80% of MSSs' energy compared to the schemes in [14], [15].

5 SCHEDULING SOLUTIONS UNDER THE MESH ARCHITECTURE

Under the mesh architecture, all SSs will organize a multihop ad hoc network for communication. To improve network throughput, some research efforts aim at reducing the *scheduling length*, which is the number of mini-slots that the BS can serve all SSs' data. Below, we introduce two scheduling solutions. One is to consider adopting spatial reuse to allow more concurrent transmissions, while another is to reduce the scheduling length (including the transmission overhead caused by guard time) to improve efficiency.

5.1 Scheduling Solution to Enhance Concurrent Transmissions

By adopting the nature of spatial reuse under the mesh architecture, the work of [17] proposes two strategies to allow more concurrent transmissions to reduce the scheduling length. Two transmissions are allowed to coexist if they do not interfere with each other. The first strategy is to try to find out those concurrent transmissions that can transmit the maximum amount of data. When all SSs have the same transmission rate, this strategy will find the maximum number of concurrent transmissions. Then,

among the concurrent transmissions being calculated, the BS selects the minimum length of bursts among these transmissions and allocates a burst with that length to each of them. The above operations are repeated until all data in each SS's queue are consumed.

Let q_i , c_i , and h_i be the queue length, transmission rate, and hop count from the BS of an SS i , respectively. The second strategy adopts six criteria to select concurrent transmissions:

- 1) The total queue length of the SSs selected for concurrent transmissions (denoted by a set \hat{S}):
 $\rho = \sum_{i \in \hat{S}} q_i$.
- 2) The total transmission rate of the SSs in \hat{S} :
 $\rho = \sum_{i \in \hat{S}} c_i$.
- 3) The queue lengths and transmission rates of the SSs in \hat{S} :
 $\rho = \sum_{i \in \hat{S}} q_i \times c_i$.
- 4) The total transmission time of the SSs in \hat{S} :
 $\rho = \sum_{i \in \hat{S}} q_i / c_i$.
- 5) The hop counts and queue lengths of the SSs in \hat{S} :
 $\rho = \sum_{i \in \hat{S}} (h_{\max} - h_i + 1) \times q_i$, where h_{\max} is the maximum hop count in the network.
- 6) The hop counts, queue lengths, and transmission rates of the SSs in \hat{S} :
 $\rho = \sum_{i \in \hat{S}} (h_{\max} - h_i + 1) \times q_i \times c_i$.

Then, for each criterion, the BS selects the set of SSs with the largest ρ value to serve. Similar to the first strategy, the BS selects the minimum length of bursts among these concurrent transmissions and allocates a burst with that length to each of them. The above operations are repeated until the queues of all SSs become empty.

By allowing more concurrent transmissions, the scheduling length is reduced so that network throughput can be improved. Considered a 5×5 grid topology with random generated traffics, the simulation results show that the sixth criterion can achieve the best performance (in terms of the scheduling length). However, since both strategies have to test all possible combinations of concurrent transmissions, the BS may encounter a high computation complexity. In addition, since the BS allocates the minimum length of bursts for all concurrent transmissions, some SSs may need to transmit their data using multiple bursts, which increases transmission overhead. These issues will be addressed in the next section.

5.2 Scheduling Solution to Reduce Scheduling Length

The work of [18] aims at regular transmissions in grid-based WiMAX mesh networks, which have been deployed in many areas such as South Africa [4]. By employing regular transmissions, not only the scheduling complexity can be reduced, but also network throughput can be improved by allowing more concurrent transmissions. The objective is to find out the optimal burst size for SSs to transmit so that the scheduling length (including the transmission overhead caused by guard time) can be minimized.

Given a grid-based WiMAX mesh network, the idea is to partition it into multiple chain-based networks and schedule each chain. Each chain has only one receiver to collect data from all other SSs along the chain. Then, the result can be extended to the whole grid-based network by letting the receiver in each chain to send data to the BS. For each chain, three possible cases may be considered:

- **There is only one source and the receiver locates at one end of the chain.** This case is the simplest one. Suppose that the interference range is fixed so that we can partition SSs into multiple disjointed groups to guarantee concurrent transmissions. In this case, the transmissions of SSs can be realized in a ‘pipeline’ manner, as shown in Fig. 6. Since all transmissions are regular, the problem is to find the optimal burst size to minimize the scheduling length. Fig. 6(a) and (b) together give an example, where SS_7 is the source with a request of four bytes. Assume that the guard time takes one mini-slot and the link rate is one byte per mini-slot. The interference range is two hops so that two SSs with a distance more than two hops can concurrently transmit their data without interfering with each other. In each cycle, three concurrent transmission flows can coexist: $SS_7 \rightarrow^{(1)} SS_6 \rightarrow^{(2)} SS_5 \rightarrow^{(3)} SS_4$, $SS_4 \rightarrow^{(1)} SS_3 \rightarrow^{(2)} SS_2 \rightarrow^{(3)} SS_1$, and $SS_1 \rightarrow^{(1)}$ receiver, where ‘ $\rightarrow^{(i)}$ ’ indicates the order of a transmission. In Fig. 6(a), the burst size is one mini-slot so that the cycle length is $[1 \text{ (guard time)} + 1 \text{ (burst size)}] \times 3 \text{ (maximum hop count in a transmission flow)} = 6$ mini-slots. Since SS_7 has four-byte data and each burst can carry one-byte data, it takes totally $4/1 = 4$ cycles for SS_7 to send all its data to SS_4 . In addition, SS_4 takes one cycle (that is, the fifth cycle) to send the last burst to SS_1 and SS_1 spends two mini-slots to forward this burst to the receiver. Therefore, the total scheduling length is $5 \text{ (the number of cycles)} \times 6 \text{ (cycle length)} + 2 \text{ (} SS_1 \text{ forwards the last burst)} = 32$ mini-slots. On the other hand, in Fig. 6(b), the burst size is two mini-slots so that each cycle takes $(1+2) \times 3 = 9$ mini-slots. Since SS_7 has four-byte data and each burst can carry two-byte data, it takes totally $4/2 = 2$ cycles to send all its data to SS_4 . In addition, SS_4 takes one cycle (that is, the third cycle) to send the last burst to SS_1 and SS_1 spends three mini-slots to forward this burst to the receiver. Therefore, the total scheduling length is $3 \times 9 + 3 = 30$ mini-slots. It can be observed that the scheduling length can be reduced if the burst size is two mini-slots. The optimal burst size can be found using the similar calculation.
- **There are multiple sources and the receiver locates at one end of the chain.** This case can be viewed as an extension of the previous case. Considering the same assumptions, Fig. 6(c) and (d) together give an example, where SS_6 and SS_3 has a request of two and four bytes, respectively. In each cycle, two concurrent transmission flows can coexist: $SS_6 \rightarrow^{(1)} SS_5 \rightarrow^{(2)} SS_4 \rightarrow^{(3)} SS_3$ and $SS_3 \rightarrow^{(1)} SS_2 \rightarrow^{(2)} SS_1 \rightarrow^{(3)}$ receiver. In Fig. 6(c), the burst size is one mini-slot so that the cycle length is $[1 \text{ (guard time)} + 1 \text{ (burst size)}] \times 3 \text{ (maximum hop count in a transmission flow)} = 6$ mini-slots. In the first two cycles, SS_6 can send all its data to SS_3 . However, SS_3 can simultaneously send its one-byte data to the receiver after the first cycle. Thus, SS_3 has to send its remaining three-byte data and forward SS_6 ’s one-byte data to the receiver, which take four extra cycles. Therefore, the total scheduling length is $[2 + 4 \text{ (the number of cycles)}] \times 6 \text{ (cycle length)} = 36$ mini-slots. On the other hand, in Fig. 6(d), the

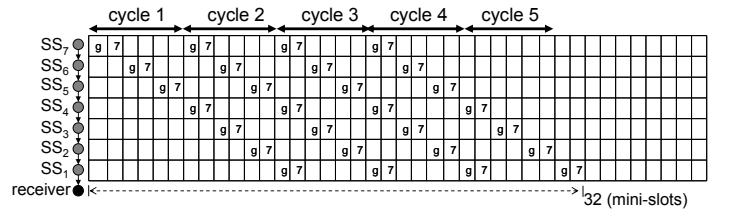
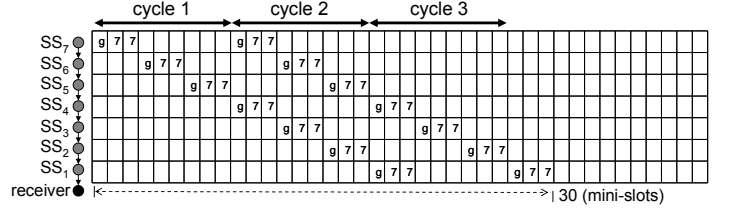
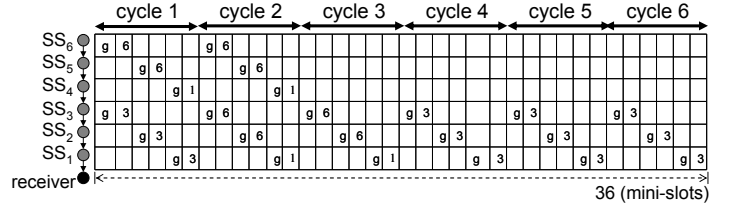
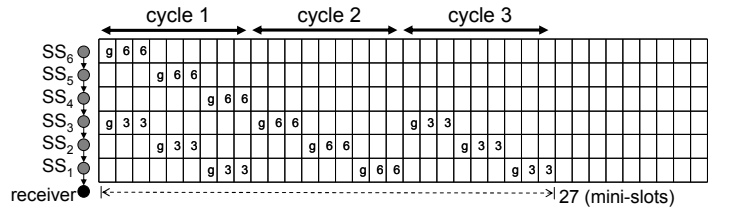
(a) One source: SS_7 with a four-byte request(b) One source: SS_7 with a four-byte request(c) Two sources: SS_6 with a two-byte request and SS_3 with a four-byte request(d) Two sources: SS_6 with a two-byte request and SS_3 with a four-byte request

Fig. 6: The case when the receiver locates at one end of the chain, where a mini-slot marked by ‘g’ is used for guard time and a mini-slot marked by a number i is used to transmit the data of SS_i . (a) The burst size is one mini-slot so that each cycle takes 6 mini-slots. The total scheduling length is 32 mini-slots. (b) The burst size is two mini-slots so that each cycle takes 9 mini-slots. The total scheduling length is 30 mini-slots. (c) The burst size is one mini-slot so that each cycle takes 6 mini-slots. The total scheduling length is 36 mini-slots. (d) The burst size is two mini-slots so that each cycle takes 9 mini-slots. The total scheduling length is 27 mini-slots.

burst size is two mini-slots so that the cycle length is $(1+2) \times 3 = 9$ mini-slots. In the first cycle, not only SS_6 can send all its data to SS_3 but also SS_3 can send its two-byte data to the receiver. Thus, SS_3 requires only two extra cycles to send its two-byte data and forward SS_6 ’s data to the receiver. Therefore, the total scheduling length is $(1+2) \times 9 = 27$ mini-slots. It can be observed that the scheduling length can be reduced if the burst size is two mini-slots. The optimal burst size can be derived following the similar calculation.

- **There are multiple sources and the receiver does not locate at either end of the chain.** In this case, the chain can be separated into a *left subchain* and a *right subchain*. We can first calculate the number of groups of SSs that are allowed to concurrent transmit:

$$k = \begin{cases} H & \text{if } 2 \leq H \leq 4 \\ 2H - 4 & \text{if } H \geq 5 \end{cases},$$

where H is the minimum hop count that two SSs can

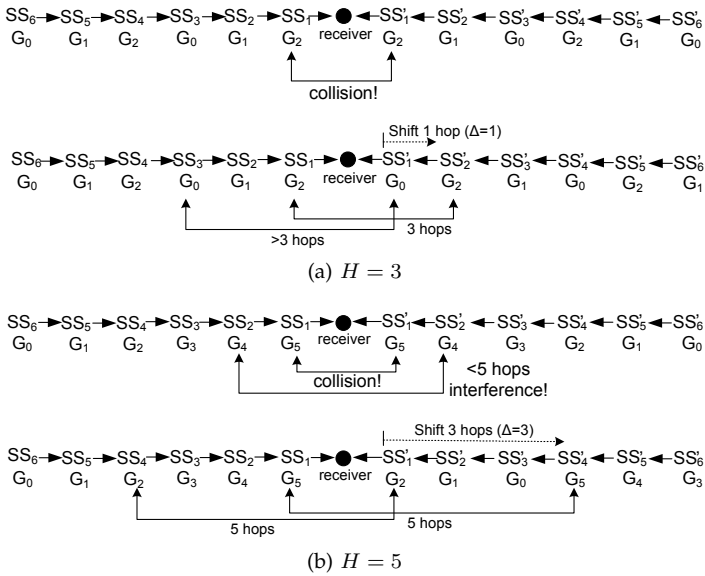


Fig. 7: The case when the receiver does not locate at either end of the chain. (a) The chain requires $k = 3$ groups of SSs and the groups in the right subchain should be shifted by $\Delta = 1$ hop. (b) The chain requires $k = 6$ groups of SSs and the groups in the right subchain should be shifted by $\Delta = 3$ hops.

concurrently transmit without interfering with each other. Then, from the end of each subchain, we first divide SSs into groups for concurrent transmission and then we ‘shift’ the groups of SSs by a number of Δ hops in the right subchain to avoid collision at the receiver, where $\Delta = 1$ if $H = 2$ and $\Delta = H - 2$ if $H \geq 3$. Fig. 7(a) and (b) together give an example. In Fig. 7(a), assuming that $H = 3$, we have $k = H = 3$ groups of SSs (that is, G_0 , G_1 , and G_2). Since both SS_1 and SS'_1 will collide at the receiver, we need to shift the groups in the right subchain by a number of $\Delta = 1$ hop. On the other hand, in Fig. 7(b), assuming that $H = 5$, we have $k = 2H - 4 = 6$ groups of SSs (that is, $G_0 \sim G_5$). Since both SS_1 and SS'_1 will collide at the receiver and SS_2 and SS'_2 will interfere with each other, we need to shift the groups in the right subchain by a number of $\Delta = H - 2 = 5$ hops. Then, we can adopt the calculation in the previous case to find out the optimal burst size. Note that the scheduling length of the whole chain will be the maximum one of both the left and right subchains.

Then, the BS adopts a *fishbone-like routing* to collect data from all SSs. In particular, the network is formed by a number of *branch chains* and one *trunk chain*, where a branch chain is a vertical chain and the trunk chain is a horizontal chain containing the BS, as shown in Fig. 8. The intersected SS of a branch chain and the trunk chain will be the receiver in that branch chain. Two branch chains parallel with a distance more than or equal to H hops are allowed to concurrently transmit. After collecting all data along each branch node, the receivers (that is, the intersected SS) will forward these data to the BS along the trunk chain. Fig. 8 give an example, where the branch chains with the same number are allowed to concurrently transmit.

By dividing a grid-based WiMAX mesh network into multiple chains and employing regular transmissions, not only the computation complexity of the scheduling solution can be reduced but also the overall guard time can be

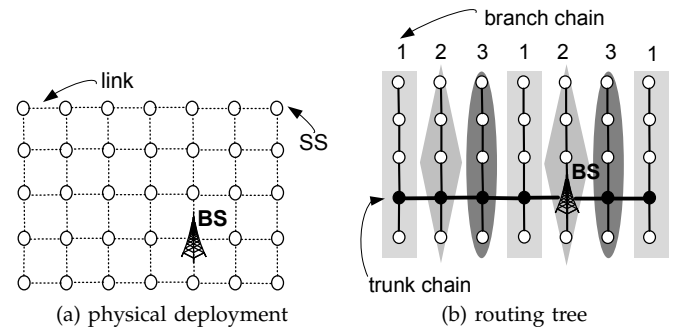


Fig. 8: The fishbone-like routing scheme in a 5×7 grid-based WiMAX mesh network, where the branch chains with the same number are allowed to concurrently transmit.

lowered down. Therefore, the scheduling length is reduced and network throughput can be improved. Considered both a chain topology with 15 nodes and a 7×7 grid topology, the simulation results show that the proposed scheme can achieve the same scheduling length but reduce up to 90% of computational complexity, compared with the scheme in [17].

6 CONCLUSION

WiMAX is developed to support broadband wireless access and its goal is to provide diverse QoS service classes under different network architectures, which makes the WiMAX scheduling problem more important and challenging. This chapter provides a comprehensive survey of recent research on scheduling solutions in WiMAX networks. Various solutions under three network architectures, including PMP, relay, and mesh, have been discussed. Table 3 gives a comparison of these scheduling solutions. For the PMP architecture, the studies of [11]–[13] adopt priority rules to support QoS of traffics. In addition, the study of [12] considers reducing burst overhead encountered in the physical layer in a low-complexity manner. The study of [13] addresses subchannel diversity and serves traffics in good subchannels first to improve network throughput. For the relay architecture, both the studies of [14]–[16] address burst allocation and adopt spatial reuse to allow concurrent transmissions in a low-complexity manner. In addition, the study of [16] tries to reduce MSSs’ transmission powers by adjusting their paths and MCSs and the groups of transmissions allowed to coexist. For the mesh architecture, both the studies of [17], [18] consider burst allocation and possibility of spatial reuse. By dividing a grid-based WiMAX mesh network into multiple chains and employing regular transmissions, the study of [18] can significantly reduce the scheduling complexity of the BS.

The IEEE 802.16 Working Group is still developing new WiMAX standards. For example, a new amendment, IEEE 802.16m [19], is designed for fourth-generation (4G) systems. The air interface of 802.16m devices is expected to achieve a transmission rate of up to 100 Mbps for mobile applications and 1 Gbps for fixed applications. The *multiple-input multiple-output (MIMO)* technique is adopted to achieve the above goal. By applying multiple antennas on both the BS and MSSs, the transmission and reception capacities can be significantly improved. In addition, to offer multiplexing rate and diversity gain, the 802.16m standard exploits the multicell MIMO to eliminate dominant intercell interferences. However, more technical issues

TABLE 3: Comparison of the features of WiMAX scheduling solutions.

reference	network architecture	QoS support	overhead reduction	subchannel diversity	burst allocation	spatial reuse	energy conservation	complexity reduction
[11]	PMP	✓						
[12]	PMP	✓	✓					✓
[13]	PMP	✓		✓				
[14], [15]	relay	✓			✓	✓		✓
[16]	relay	✓			✓	✓	✓	✓
[17]	mesh	✓			✓	✓		
[18]	mesh	✓	✓		✓	✓		✓

need to be addressed, such as the synchronization issue, how to report the channel state information, and how to get the pre-coding information from the backhaul network. These issues make the WiMAX scheduling problem more challenging and interesting.

REFERENCES

- [1] C. Eklund, R.B. Marks, K.L. Stanwood, and S. Wang, "IEEE standard 802.16: A technical overview of the wirelessMAN air interface for broadband wireless access," *IEEE Comm. Magazine*, vol. 40, no. 6, pp. 97–107, 2002.
- [2] B. Li, Y. Qin, C.P. Low, and C.L. Gwee, "A survey on mobile WiMAX," *IEEE Comm. Magazine*, vol. 45, no. 12, pp. 70–75, 2007.
- [3] WiMAX deployment in Asia. [Online]. Available: <http://www.goingwimax.com/wimax-deployment-in-asia-4192/>
- [4] D. Johnson, "Evaluation of a single radio rural mesh network in South Africa," *Proc. IEEE Int'l Conf. Information and Comm. Technologies and Development*, 2007.
- [5] IEEE Standard 802.16d-2004, "IEEE standard for local and metropolitan area networks part 16: air interface for fixed broadband wireless access systems," 2004.
- [6] IEEE Standard 802.16e-2005, "IEEE standard for local and metropolitan area networks part 16: air interface for fixed and mobile broadband wireless access systems amendment 2: Physical and medium access control layers for combined fixed and mobile operation in licensed bands and corrigendum 1," 2006.
- [7] IEEE Standard 802.16j-2009, "IEEE standard for local and metropolitan area networks part 16: air interface for broadband wireless access systems amendment 1: Multiple relay specification," 2009.
- [8] Y. Ben-Shimol, I. Kitroser, and Y. Dinitz, "Two-dimensional mapping for wireless OFDMA systems," *IEEE Trans. Broadcasting*, vol. 52, no. 3, pp. 388–396, 2006.
- [9] T. Ohseki, M. Morita, and T. Inoue, "Burst construction and packet mapping scheme for OFDMA Downlinks in IEEE 802.16 systems," *Proc. IEEE Global Telecomm. Conf.*, pp. 4307–4311, 2007.
- [10] T. Wang, H. Feng, and B. Hu, "Two-dimensional resource allocation for OFDMA system," *Proc. IEEE Int'l Conf. Comm. Workshops*, 2008.
- [11] Q. Liu, X. Wang, and G.B. Giannakis, "A cross-layer scheduling algorithm with QoS support in wireless networks," *IEEE Trans. Vehicular Technology*, vol. 55, no.3, pp. 839–847, 2006.
- [12] J.M. Liang, J.J. Chen, Y.C. Wang, Y.C. Tseng, and B.S.P. Lin, "Priority-Based scheduling algorithm for downlink traffics in IEEE 802.16 networks," *Proc. IEEE Asia Pacific Wireless Comm. Symp.*, 2009.
- [13] X. Zhu, J. Huo, X. Xu, C. Xu, and W. Ding, "QoS-guaranteed scheduling and resource allocation algorithm for IEEE 802.16 OFDMA system," *Proc. IEEE Int'l Conf. Comm.*, pp. 3463–3468, 2008.
- [14] V. Genc, S. Murphy, and J. Murphy, "An interference-aware analytical model for performance analysis of transparent mode 802.16j systems," *Proc. IEEE Global Comm. Conf.*, 2008.
- [15] V. Genc, S. Murphy, and J. Murphy, "Analysis of transparent mode IEEE 802.16j system performance with varying numbers of relays and associated transmit power," *Proc. IEEE Wireless Comm. and Networking Conf.*, 2009.
- [16] J.M. Liang, Y.C. Wang, J.J. Chen, J.H. Liu, and Y.C. Tseng, "Efficient resource allocation for energy conservation in uplink transmissions of IEEE 802.16j transparent relay networks," *Proc. ACM Int'l Conf. Modeling, Analysis and Simulation of Wireless and Mobile Systems*, pp. 187–194, 2010.
- [17] J. Zhang, H. Hu, L. Rong, and H. Chen, "Cross-Layer scheduling algorithms for IEEE 802.16 based wireless mesh networks," *Wireless Personal Comm.*, vol. 51, no. 3, pp. 615–634, 2009.
- [18] J.M. Liang, H.C. Wu, J.J. Chen, and Y.C. Tseng, "Mini-slot scheduling for IEEE 802.16d chain and grid mesh networks," *Computer Comm.*, vol. 33, no. 17, pp. 2048–2056, 2010.
- [19] IEEE Standard 802.16m/D11, "IEEE draft amendment standard for local and metropolitan area networks part 16: air interface for broadband wireless access systems—advanced air interface," 2011.