



# DISSERTATION

In Partial Fulfillment of the Requirements  
for the Degree of Doctor of Philosophy

from TELECOM ParisTech

Specialization: Computer Science and Networking

**Ikbal Chammakhi Msadaa**

## **QoS Management and Performance Analysis of WiMAX Networks in Fixed and Highly Mobile Environments**

Defense scheduled on the 6th of October 2010 before a committee composed of:

Reporters	André-Luc Beylot, IRIT-ENSEEIH, France Thierry Turletti, INRIA Sophia-Antipolis, France
Examiners	Christian Bonnet, EURECOM, France Guy Juanole, LAAS-CNRS Toulouse, France
Thesis supervisor	Fethi Filali, QU Wireless Innovations Center, Qatar





# THÈSE

présentée pour obtenir le grade de  
Docteur de TELECOM ParisTech

Spécialité : Informatique et Réseaux

**Ikbal Chammakhi Msadaa**

## **Gestion de la QoS et Evaluation des Performances des Réseaux WiMAX dans Les Environnements Fixes et à Forte Mobilité**

Soutenance prévue pour le 6 octobre 2010 devant le jury composé de:

Rapporteurs	André-Luc Beylot, IRIT-ENSEEIH, France Thierry Turletti, INRIA Sophia-Antipolis, France
Examineurs	Christian Bonnet, EURECOM, France Guy Juanole, LAAS-CNRS Toulouse, France
Directeur de thèse	Fethi Filali, QU Wireless Innovations Center, Qatar





# Abstract

Driven by the growing demand for high-speed broadband wireless services, Worldwide Interoperability for Microwave Access (WiMAX) technology has emerged as a competitive alternative to wireline broadband access solution. WiMAX technology, considered in this thesis, offers an IP-based framework that provides high data rates at medium and long range with the ability of supporting fixed, nomadic, portable, and mobile access. Moreover, based on the IEEE 802.16 standard, the technology provides a set of built-in QoS mechanisms to support heterogeneous classes of traffic including data, voice and video. The IEEE 802.16 standard, however, leaves unstandardized the resource management and scheduling mechanisms, which are crucial components to guarantee QoS performance for these services.

In this thesis, we evaluate the performance of IEEE 802.16 based WiMAX technology in both fixed and highly mobile environments. More particularly, Mobile WiMAX is investigated as a vehicular-to-infrastructure (V2I) communication medium since it is expected to play a major role in intelligent transportation systems. The technology is indeed the only mobile broadband technology currently in use.

Moreover, we address in this thesis most of the resource management and scheduling issues that have been left open with the objective of defining an architecture that fulfills the QoS expectations of the five classes of applications addressed by the IEEE 802.16 standard. In fact, after surveying, classifying and comparing different scheduling and admission control mechanisms proposed in this work-in-progress area, we propose two QoS solutions. Both solutions address point-to-multipoint (PMP) 802.16 systems operating in time division duplex mode (TDD) mode.

The first solution includes a hierarchical scheduling algorithm that adapts the DL/UL allocations on a frame-by-frame basis to serve unbalanced traffic. The amounts of these bandwidth grants are set by the connection admission control (CAC) module that adopts a Max-Min fairness approach making efficient and fair use of the available resources. The proposed solution takes into account the link adaptation capability supported by WiMAX and the data rate constraints of the different types of services.

The second QoS solution presented in this thesis is a multi-Constraints Scheduling Strategy (mCoSS) that is designed for both OFDM or band-AMC OFDMA air interfaces. Unlike the first QoS solution, mCoSS supposes the use of a predefined DL/UL ratio set by the operator. In addition to data rate constraints, mCoSS offers the advantage (compared to the first solution) of supporting delay constraints of real-time applications and handling bursty traffics. mCoSS is based on a modified dual-bucket traffic shaping mechanism configured on a per-flow basis. This shaping mechanism is combined with a two-rounds scheduling strategy which reflects (i) at the first round, the minimum data rates and latency requirements the BS or MS is committed to provide and (ii) at the second round, the efficiency and fairness of the resources management since the remaining bandwidth is shared in this round using a simple weighted fair queuing (WFQ) strategy; Nevertheless, the allocations should remain within the thresholds set by the dual-bucket shaping mechanism.

---



## Résumé

Durant les deux dernières décennies, le développement dans le domaine des réseaux télécoms a façonné notre quotidien grâce au succès de l'accès sans fil et a créé chez le grand public un besoin accru en débit. Les utilisateurs souhaitent en effet avoir une qualité de service équivalente à celle perçue dans les réseaux filaires.

Entre autre solutions candidates, la technologie WiMAX (Worldwide Interoperability for Microwave Access), à laquelle nous nous intéressons dans le cadre de cette thèse, a émergé afin de répondre à ces nouveaux besoins. En plus d'être sans fil et à haut débit, la technologie WiMAX est basée sur IP et mobile. D'ailleurs, ces fonctionnalités la positionne comme une technologie de pointe qui vient à bout des tarifs élevés des technologies 3G et de la mobilité limitée du WiFi. En outre, le WiMAX Mobile est une réalité et est en train d'être déployé aux Etats Unis, au Japon, en Corée, en Europe, en Australie et un peu partout dans le monde. C'est en fait la seule technologie haut débit mobile en cours d'utilisation. Il y a même des discussions en cours concernant l'éventuelle sélection du WiMAX comme standard International Mobile Telecommunications (IMT)-advanced.

Le WiMAX est basé sur la famille de standards et amendements IEEE 802.16 spécifiant les couches MAC et PHY pour l'accès fixe, nomade, portable et mobile. De plus, la technologie présente un ensemble d'éléments clefs: (1) l'utilisation de l'orthogonal frequency division multiplex (OFDM), (2) le duplexage temporel et fréquentiel (TDD et FDD), (3) le support de la modulation et du codage adaptatif (AMC) et (4) des techniques d'antennes avancées telles que les antennes multiple input, multiple output (MIMO), (5) une sécurité robuste et (6) des éléments permettant de supporter les besoins en qualité de service (QoS) de plusieurs types de trafics. Dans le cadre de cette thèse, nous nous intéressons justement à la gestion de la qualité de service dans les réseaux WiMAX et plus particulièrement aux problèmes d'ordonnancement et de contrôle d'admission (CAC) qui en découlent. En effet, bien qu'il présente des éléments permettant de véhiculer des données, de la voix ainsi que de la vidéo, le protocole MAC de la norme IEEE 802.16 laisse ouverts les problèmes rattachés à l'ordonnancement et au contrôle d'admission; des éléments cruciaux pour l'amélioration de la QoS perçue par les utilisateurs. Dans cette thèse<sup>1</sup>, nous évaluons les performances des réseaux WiMAX dans les environnements fixes et à forte mobilité. Nous étudions plus particulièrement le potentiel et les limites de l'utilisation du WiMAX Mobile en tant que médium de communications véhicule-à-infrastructure (V2I). Nous attaquons, dans le cadre de cette thèse, essentiellement aux problèmes de gestion de ressources laissés ouverts par le standard IEEE 802.16. Le reste du manuscrit est organisé comme décrit dans la section qui suit.

---

<sup>1</sup>Ce travail a été soutenu par le projet WiNEM (WiMAX Network Engineering and Multihoming) sous la subvention No. 2006 TCOM005 05 et par les membres industriels d'EURECOM: BMW Group, Cisco, Monaco Telecom, Orange, SAP, SFR, Sharp, STEricsson, Swisscom, Symantec, et Thales.

---

## Structure et Contributions

### Chapitre 1: Un Aperçu de la Technologie WiMAX

L'objectif de ce premier chapitre est de donner un aperçu général de la technologie WiMAX. Nous commençons donc par passer en revue le processus de standardisation de la famille de standards IEEE 802.16. Puis, nous présentons les différentes interfaces physiques et bandes de fréquences correspondantes. Ensuite, la couche physique est décrite en accordant un intérêt plus particulier à la technique de modulation et de codage adaptatif (AMC) supportée par la technologie WiMAX. La couche MAC est également décrite mais d'une manière plus brève; seules les fonctionnalités de base, nécessaires à la compréhension de l'étude de performance menée dans le Chapitre 2, sont présentées au niveau de ce chapitre. Tous les concepts relatifs au support de la qualité de service (QoS) sont détaillés dans le Chapitre 3. En effet, étant donné le nombre important de concepts introduits par le standard IEEE 802.16 à cet effet, nous avons préféré leur dédier un chapitre en entier.

### Chapitre 2: Analyse de Performances des Réseaux WiMAX basés sur OFDM

Dans ce chapitre, nous évaluons les performances théoriques maximales des systèmes WiMAX. Le débit de saturation qui pourrait être atteint dans des réseaux WiMAX est calculé à travers plusieurs scénarios où l'on fait varier par exemple la durée de la trame physique, la bande passante des canaux ou le schéma de modulation et de codage (MCS). Un modèle analytique a été développé en se basant sur des propriétés techniques et des profils systèmes spécifiés par le standard IEEE 802.16 pour des systèmes utilisant l'interface physique WirelessMAN-OFDM.

Certaines parties de ce chapitre ont été publiées dans:

- Ikbal Chammakhi Msadaa and Fethi Filali. On the Performance Bounds of OFDM-based 802.16 Broadband Wireless Networks. In WCNC 2008, IEEE Wireless Communications and Networking Conference, Apr. 2008.

### Chapitre 3: Support de la QoS dans les Réseaux WiMAX

Le standard IEEE 802.16 définit un protocole MAC orienté connexion qui est conçu pour s'adapter à des applications avec des besoins divers en QoS. Néanmoins, plusieurs problèmes, rattachés notamment à la gestion de ressources, avaient été laissés ouverts. L'objectif principal de ce chapitre est de fournir une vision plus claire de ce qui est supporté ou non afin d'améliorer la QoS perçue par les utilisateurs dans les réseaux WiMAX. Pour ce faire, nous commençons par décrire les principaux éléments mis en place par le standard afin de répondre aux besoins de trafics hétérogènes. Ensuite, nous proposons une architecture générique qui incorpore les principaux composants nécessaires à la mise en place d'une politique de gestion de la QoS dans les systèmes WiMAX. La dernière section de ce chapitre est consacrée aux problèmes d'ordonnancement et de contrôle d'admission. Plus précisément, nous mettons en évidence les principaux défis à relever lors de la conception d'une solution d'ordonnancement et/ou de contrôle d'admission (CAC) pour les réseaux WiMAX.

Certaines parties de ce chapitre ont été publiées dans:

---

- Ikbal Chammakhi Msadaa, Fethi Filali, and Farouk Kamoun. An 802.16 Model for NS2 Simulator with an Integrated QoS architecture. In SIMUTools' 08, 1st International Conference on Simulation Tools and Techniques for Communications, Networks and Systems, Mars 2008.

#### **Chapitre 4: Ordonnement et CAC: Etude et Taxonomie**

De nombreux chercheurs ont été concernés par les problèmes d'ordonnement et de contrôle d'admission dans les réseaux WiMAX. Dans ce chapitre, nous faisons un état de l'art des travaux existant dans la littérature, classons et analysons les solutions proposées dans ce domaine. Certaines parties de ce chapitre ont été publiées dans:

- Ikbal Chammakhi Msadaa, Daniel Câmara, and Fethi Filali. Scheduling and CAC in IEEE 802.16 fixed BWNs : a Comprehensive Survey and Taxonomy. "IEEE Communications Surveys & Tutorials", 12(4):459–487, 2010.
- Tijani Chahed, Ikbal Chammakhi Msadaa, Rachid Elazouzi, Fethi Filali, Salah-Eddine Elayoubi, Benoit Fourestié, Thierry Peyre, and Chadi Tarhini. WiMAX Network Capacity and Radio Resource Management. Book chapter in "Radio Resources Management in WiMAX : From theoretical capacity to system simulations", ISBN: 9781848210691, Feb. 2009.

#### **Chapitre 5: Ordonnement Adaptatif et Contrôle d'Admission Max-Min**

Bien qu'incluant des éléments qui permettent de supporter la QoS, le protocole MAC 802.16 ne constitue pas une solution complète qui puisse répondre aux besoins de diverses applications. En effet, les problèmes d'ordonnement et de gestion de ressources ont été laissés ouverts. Dans ce chapitre, nous proposons une nouvelle architecture de QoS pour les systèmes WiMAX point-à-multipoints (PMP) opérant en mode TDD et utilisant l'interface physique WirelessMAN-OFDM. Cette architecture inclut une politique de contrôle d'admission et un algorithme d'ordonnement hiérarchique. La solution CAC adopte un schéma d'équité Min-Max utilisant d'une manière efficace et équitable les ressources disponibles. L'algorithme d'ordonnement proposé ajuste d'une manière flexible la bande passante entre le lien descendant et le lien ascendant s'adaptant ainsi à un éventuel trafic asymétrique. Cette façon d'opérer prend en considération la technique de modulation et de codage adaptatif mise en oeuvre par les systèmes WiMAX ainsi que les contraintes de débits de chaque connexion. La solution proposée se révèle, à travers les simulations, efficace et capable de s'adapter aux besoins en débits des divers types de services visés par le standard IEEE 802.16.

Certaines parties de ce chapitre ont été publiées dans:

- Ikbal Chammakhi Msadaa, Fethi Filali, and Farouk Kamoun. An adaptive QoS Architecture for IEEE 802.16 Wireless Broadband Networks. In MASS 2007, 4th IEEE International Conference on Mobile Ad-hoc and Sensor Systems, Oct. 2007.
-

## Chapitre 6: mCoSS: une Stratégie d'Ordonnement multi-Contraintes pour les Réseaux WiMAX

Nous proposons dans ce chapitre une stratégie d'ordonnement multi-contraintes baptisée "multi-Constraints Scheduling Strategy" (mCoSS) qui maximise le niveau de QoS aussi bien pour les applications temps réel que pour celles tolérantes aux délais. mCoSS s'attaque à des contraintes qui n'avaient pas été considérées dans la solution décrite dans le Chapitre 5 à savoir la sporadicité et les besoins en délais des trafics temps réel. Selon cette stratégie, l'accès au réseau est régulé par un shaper inspiré du mécanisme de double seau à jetons qui permet d'avoir un trafic sporadique tout en protégeant les connexions conformes au contrat de service de celles qui sont gourmandes en bande passante. Cette version modifiée du double seau à jetons est combinée à un algorithme d'ordonnement à deux étapes réfléchissant les deux niveaux de service attendus par une connexion donnée. Dans une première étape, le débit minimum réservé ainsi que les contraintes de délai sont assurés. La deuxième étape consiste à répartir équitablement le reste de bande passante entre les différents flux en utilisant la politique weighted fair queuing (WFQ). La politique de demande de bande passante adoptée dans cette stratégie profite de la multitude de techniques proposées par le standard IEEE 802.16e et adapte le choix de la technique la plus appropriée aux contraintes de QoS des flux ainsi qu'à la disponibilité des ressources radio. D'autres contraintes telles que l'AMC et la protection des trafics BE de la famine sont également considérées dans la stratégie proposée.

## Chapitre 7: WiMAX Mobile: un Médium de Communications V2I

Le forum WiMAX estime que la technologie WiMAX serait déployée en majorité dans sa version mobile. Et qui dit mobilité, dit hétérogénéité de réseaux. De ce fait, nous nous attaquons à la technologie WiMAX dans le contexte mobile et hétérogène des systèmes de transport intelligents (ITS). Ces systèmes ont fait l'objet depuis les années 80 d'une stratégie mondiale qui vise à résoudre plusieurs de nos soucis de transport quotidiens. Ces systèmes permettraient en effet aux gens d'atteindre leurs destination d'une manière sûre, efficace et confortable. Afin d'atteindre ces objectifs, plusieurs technologies d'accès radio (RAT) telles que l'UMTS, le WiMAX ou encore la technologie 5.9 GHz ont été proposées pour la nouvelle génération de systèmes de transport intelligents.

En plus de la technologie 5.9 GHz, qui est spécialement dédiée aux réseaux véhiculaires, le WiMAX mobile est attendu comme une technologie qui jouerait un rôle important dans les ITS étant donné que c'est la seule technologie haut débit mobile en cours d'utilisation.

Dans ce chapitre, nous comparons le WiMAX mobile (basé sur le standard IEEE 802.16e) à la technologie 5.9 GHz (basée sur l'imminent standard IEEE 802.11p). Nous étudions, par simulation, le potentiel et les limites des deux technologies en tant que média de communications véhicule-à-infrastructure (V2I). Les performances des deux systèmes sont évaluées pour différentes vitesses de véhicule, différents débits et différents déploiements de réseaux.

Certaines parties de ce chapitre ont été publiées dans:

- Ikbal Chammakhi Msadaa, Pasquale Cataldi, and Fethi Filali. A Comparative Study between 802.11p and Mobile WiMAX-based V2I Communication Networks. In NGMAST

---

2010, 4th International Conference on Next Generation Mobile Applications, Services and Technologies, July 2010.

## **Annexe A: Sujets Relatifs à la Gestion de la Mobilité dans les réseaux WiMAX**

Le forum WiMAX estime que plus de 133 millions de personnes utiliseraient la technologie WiMAX d'ici 2012. Parmi ces utilisateurs potentiels, plus de 70% utiliseraient l'implémentation mobile de cette technologie. De ce fait, la gestion de la mobilité constitue un challenge de taille pour ces 70% d'utilisateurs WiMAX.

Cette annexe est consacrée à cette problématique. Elle décrit en effet les concepts et mécanismes introduits par le standard IEEE 802.16e—l'amendement du standard IEEE 802.16d-2004—qui apporte des améliorations qui concernent surtout la gestion de la mobilité. Nous couvrons également, à travers cette annexe, les principaux sujets relatifs à la mobilité dans les réseaux WiMAX et mettons en évidence les sujets de recherche qui sont encore ouverts aux contributions.

Certaines parties de cette annexe ont été publiées dans:

- Ikbal Chammakhi Msadaa, Daniel Câmara, and Fethi Filali. Mobility Management in WiMAX Networks. Book chapter in "WiMAX Security and Quality of Service : An End-to-End Perspective". ISBN : 978-0-470-72197-1. Seok-Yee Tang and Peter Muller and Hamid Sharif Ed., July 2010.

Dans ce résumé, nous développons certaines de ces contributions.

## **Analyse de Performances des Réseaux WiMAX basés sur OFDM**

Dans cette thèse, un modèle analytique original est développé afin d'étudier les performances théoriques maximales des systèmes 802.16 basés sur OFDM. Ce modèle analytique est développé conformément aux spécifications du standard IEEE 802.16 [1]. En se basant sur cette étude, plusieurs scénarios ont été considérés afin d'évaluer les performances théoriques maximales des systèmes WiMAX sous différentes configurations des paramètres MAC et PHY. Les résultats obtenus mettent en évidence l'importance de considérer l'overhead MAC et PHY lors de l'évaluation de performances des systèmes IEEE 802.16. En effet cet overhead, qui est souvent ignoré ou grossièrement estimé dans des travaux de recherches, pourrait constituer entre 40 et 90 % de la totalité de la trame, en fonction de la taille des PDUs et des profils systèmes considérés. Aussi avons nous montré à travers cette étude analytique que l'utilisation d'une bande passante plus large n'implique pas forcément une amélioration conséquente des performances au niveau MAC. En examinant l'effet de la fragmentation et de l'agrégation sur ces performances, nous démontrons également que celle-ci pourrait nettement améliorer les débits obtenus notamment dans le cas de trafics transportant des paquets de taille fixe.

## **Support de la QoS dans les Réseaux WiMAX**

Le standard IEEE 802.16 définit un protocole MAC orienté connexion. Chaque connexion est associée à un service flow (SF) caractérisé par un ensemble de paramètres de QoS reflétant les contraintes en débit et/ou délai de l'application correspondant à ce flux. Le Tableau 1 dresse la liste des paramètres de QoS à spécifier lors de la création d'une nouvelle connexion correspondant

---

Traffic/Applications Characteristics	real-time, fixed-rate data, Fixed/Variable length PDUs		real-time, variable bit rates, requiring guaranteed data rate and delay		real-time, variable bit rates, requiring guaranteed data rate and delay		requiring guaranteed data rate, insensitive to delays		No rate or delay requirement	
	DL	UL	DL	UL	DL	UL	DL	UL	DL	UL
Downlink (DL)/ Uplink (UL)										
Maximum Sustained Traffic Rate	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Minimum Reserved Traffic Rate	✓	✓	✓	✓	✓	✓	✓	✓	—	—
Maximum Latency	✓	✓	✓	✓	✓	✓	—	—	—	—
Tolerated Jitter	✓	✓	✓	✓	—	—	—	—	—	—
Request/Transmission Policy	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Traffic Priority	—	—	✓	✓	✓	✓	✓	✓	—	—
Request/Grant Scheduling Type	—	✓ (UGS)	—	✓ (ertPS)	—	✓ (rtPS)	—	✓ (nrtPS)	—	✓ (BE)
Unsolicited Grant Interval	—	✓	—	✓	—	—	—	—	—	—
Unsolicited Polling Interval	—	—	—	—	—	✓	—	—	—	—
SDU Size (If fixed length SDU)	✓	✓	—	—	—	—	—	—	—	—
Example of application	T1/E1, VoIP without VAD		VoIP with VAD		MPEG video		FTP		HTTP, SMTP	

Table 1: Mandatory QoS parameters for each scheduling service

à telle ou telle catégories d'applications. En outre, pour les connexions du lien ascendant, le standard IEEE 802.16 définit cinq "request/grant scheduling types", à savoir:

- unsolicited grant service (UGS),
- extended real-time polling service (ertPS),
- real-time polling service (rtPS),
- non-real-time polling service (nrtPS),
- et best effort (BE).

Il est à noter que les paramètres de qualité de service sont les mêmes pour un type d'application donné que celle-ci soit sur le DL ou le UL et pourtant les "request/grant scheduling types" ne sont associées qu'aux connexions UL. Les noms de ces types reflètent en fait la manière dont la bande passante est demandée ou allouée par la MS et la BS, respectivement pour les connexions du lien ascendant. D'ailleurs, le standard propose une multitude de techniques à cet effet.

Ce qu'il faudrait toutefois retenir est que quelque soit le mécanisme de demande et d'allocation de la bande passante, celle-ci est toujours demandée par flux mais accordée par MS. En d'autres termes la station de base répond aux besoins de plusieurs connexions UL d'une même MS sous forme d'une allocation agrégée et c'est à la MS de décider de la manière dont ces ressources seraient réparties puisqu'elle possède une perception plus accrue et plus à jour des différents besoins.

D'ailleurs ceci nous ramène à une architecture d'ordonnancement à trois composantes majeures: deux du côté de la BS (l'ordonnanceur DL et l'ordonnanceur UL) et une du côté de la MS pour les connexions UL. De plus, comme le standard ne définit pas la manière dont ces composantes interagissent entre elles ni la manière dont les différents concepts introduits pour gérer la QoS pourraient être réunis au sein d'une même architecture, nous proposons dans cette thèse une architecture de QoS qui répond à cette problématique. L'architecture que nous proposons et qui est illustrée par la Figure 1 se veut d'être un cadre assez générique qui pourrait servir de base pour concevoir des solutions d'ordonnancement et de contrôle d'admission pour les réseaux WiMAX.

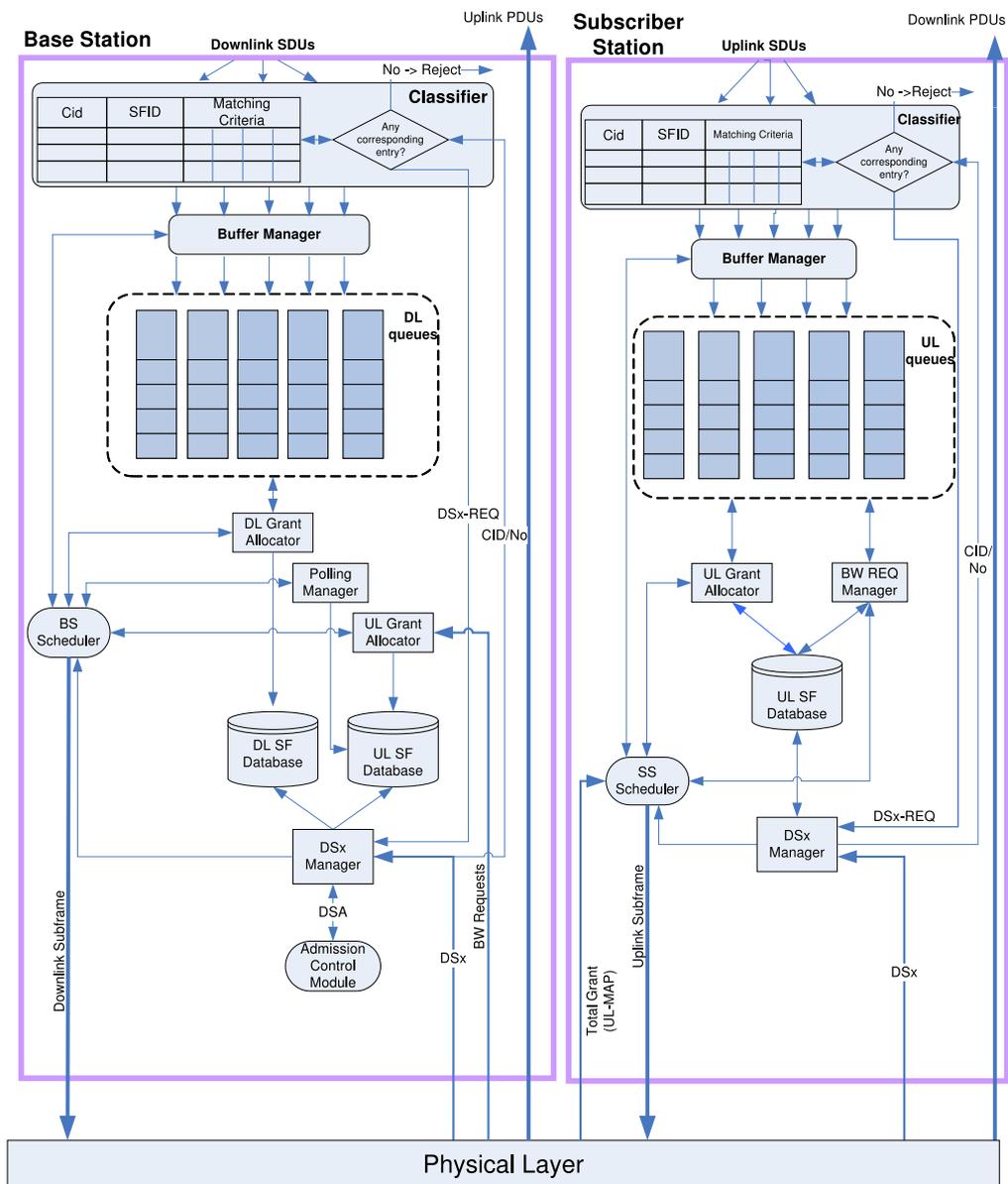


Figure 1: QoS architecture Design

## Ordonnancement et CAC: Etude et Taxonomie

Tel que le montre la Figure 2, les approches adoptées dans la littérature en concevant une solution d'ordonnancement pourraient être divisées en trois principales catégories.

1. La première est basée sur une stratégie de gestion de files d'attente où les auteurs traitent le problème comme tel et essaient de trouver la discipline de gestion de file d'attente la plus appropriée et qui pourrait au mieux répondre aux contraintes de QoS des différents types de services visés par le standard IEEE 802.16 [1, 2]. Dans cette première catégorie, deux types de structures reviennent assez souvent: soit une structure simple consistant en général en une seule politique de gestion de files d'attente appliquée à toutes les catégories d'applications [3, 4, 5] soit alors une structure hiérarchique plus élaborée, comme dans [6, 7, 8, 9, 10, 11, 12, 13, 14], basée sur deux ou plusieurs niveaux d'ordonnements reflétant les différents niveaux des décisions d'ordonnements prises.
2. Une seconde catégorie où le problème d'ordonnement est formulé sous forme d'un problème d'optimisation dont l'objectif est de maximiser les performances du système sujet à des contraintes reflétant en général les contraintes de QoS des différentes classes de services [15, 16, 17, 18, 19, 20, 21, 22, 23].
3. La troisième catégorie qui pourrait être rencontrée dans la littérature est fondée sur une approche cross-layer basée en général sur une architecture cross-layer. L'objectif de cette architecture est d'optimiser la communication entre deux [24, 25, 26, 27, 28] voire trois [29, 30] couches réseau et ainsi améliorer les performances du système.

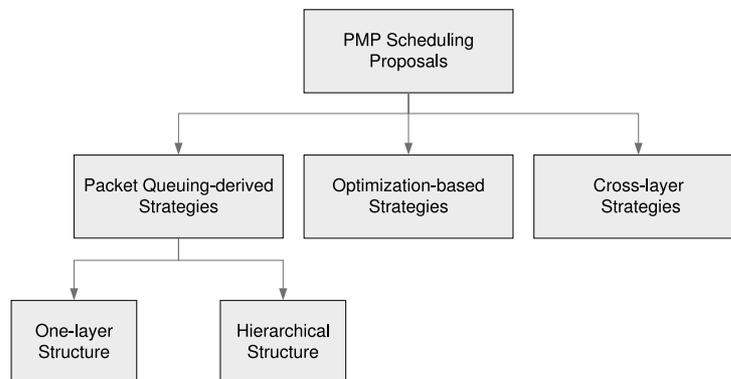


Figure 2: Classification of the scheduling strategies of IEEE 802.16 PMP mode

## mCoSS: une Stratégie d'Ordonnancement multi-Contraintes pour les Réseaux WiMAX

Nous tentons, à travers cette thèse de rassembler les différentes pièces du puzzle "gestion de ressources dans les réseaux WiMAX" en s'attaquant aux problèmes qui sont toujours ouverts. Dans cette perspective, nous proposons une stratégie d'ordonnancement multi-contraintes baptisée mCoSS (multi-Constraints Scheduling Strategy) qui définit les opérations d'ordonnement aussi bien coté BSs que coté MSs. La dite stratégie est décrite à travers un ensemble d'algorithmes

qui maximisent le degré de satisfaction en QoS des trafics temps-réel et ceux tolérants aux délais en terme de débit et de délai. mCoSS pourrait très bien s'appliquer à des environnement OFDM que band-AMC OFDMA.

L'accès au canal est régulé par le moyen d'un shaper inspiré du mécanisme de double seau à jetons qui rend possible la sporadicité tout en protégeant les trafics qui se conforment au contrat de service des trafics gourmands en bande passante. Ce mécanisme de double seau à jetons est combiné à un algorithme d'ordonnancement opérant en deux étapes. Dans un premier temps, le débit minimum réservé ainsi que les contraintes de délai sont satisfaits. Puis l'équité entre les différents flux est assurée grâce à l'utilisation de la politique d'ordonnancement weighted fair queuing (WFQ) pour partager le reste de la bande passante. La politique de demande de la bande passante profite de la multitude de techniques proposées par le standard IEEE 802.16 en adaptant le choix de la technique la plus appropriée en fonction de la quantité de ressources disponibles et des contraintes QoS du flux considéré. D'autres problèmes tels que la famine des trafics BE et la mise en oeuvre de la modulation et codage adaptatif sont également pris en considération dans notre stratégie d'ordonnancement. Afin d'évaluer cette solution, nous l'avons implémentée dans le simulateur Qualnet et l'avons comparée aux disciplines Strict Priority (SP) et une variante du WFQ. Les résultats obtenus montre un compromis intéressant entre équité et efficacité avec un respect des contraintes de qualité de service des différentes connexions.

Dans ce qui suit, nous commençons par expliquer l'idée du shaper de trafic basé sur une version modifiée du double seau à jetons. Ensuite, nous détaillons les algorithmes d'ordonnancement à deux étapes pour finir avec une évaluation de performances de la stratégie proposée.

### Une Version Modifiée du Mécanisme de Double Seau à Jetons

Afin d'assurer une QoS pour divers types de trafics, il est important d'implémenter un mécanisme de shaping (ou lissage) afin de contrôler le volume de trafic entrant en réseau et isoler ainsi les trafics gourmands en bande passante. Les deux mécanismes de lissage les plus répandus en ingénierie de trafic sont: le "leaky bucket" (ou seau percé) et le "token bucket" (seau à jetons). Le seau percé constitue un mécanisme à travers lequel un flux est lissé de manière à être transmis dans le réseau à un débit constant. Le seau à jetons quant à lui, tout en assurant un contrôle du débit, permet une certaine sporadicité limitée par un seuil configurable.

Afin de répondre aux besoins de certaines catégories d'applications visées par le WiMAX, nous choisissons la deuxième alternative (à savoir le seau à jetons) pour modéliser notre shaper. Plus particulièrement, nous utilisons la variante seaux à jetons multiples. Nous associons chaque flux  $i$  à deux seaux correspondant au débit minimum réservé  $R_{min}^i$  et au débit maximum soutenu  $R_{max}^i$ . Ces doubles seaux reflètent en fait les limites inférieures et supérieures du service à fournir à un flux donné. Chaque seau est défini à travers trois paramètres: la taille du burst (ou rafale), le débit moyen et l'intervalle de temps. La Figure 3 représente la structure en double seau associée à chaque flux de service. Le premier seau est caractérisé par:

- un débit moyen, appelé aussi "committed information rate" ( $CIR$ ), qui spécifie la quantité de données qui pourrait être transmise en moyenne par unité de temps.
- un intervalle de temps  $T_c$ , appelé aussi intervalle de mesure; il spécifie le quantum de temps en seconde par rafale.
- la taille du burst/rafale, appelé également "committed burst size" ( $B_c$ ); elle correspond à la quantité de trafic qui pourrait être transmise par burst durant un intervalle de mesure donné.

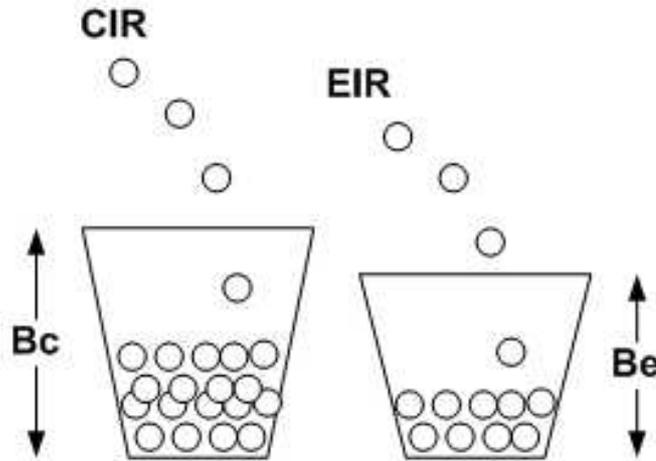


Figure 3: A Dual-Bucket Shaping Mechanism

Les trois paramètres sont reliés comme suit:  $CIR = \frac{B_c}{T_c}$ . Nous fixons  $CIR$  au débit minimum réservé  $R_{min}^i$ , et  $T_c$  à l'intervalle d'allocation  $I_{gr}^i$  caractérisant le flux  $i$ . Pour un trafic temps réel, ce paramètre correspond au délai maximum  $L_{max}^i$ . Pour les flux tolérants aux délais, ce paramètre ne devrait pas dépasser l'intervalle de polling (pour nrtPS) et pourrait être fixé en fonction de l'intervalle moyen de transmission du flux. L'introduction de ce paramètre est nécessaire afin de définir la fréquence à laquelle les allocations sont faites pour chaque flux. En effet, le standard ne spécifie pas l'intervalle sur lequel les moyennes  $R_{min}^i$  et  $R_{max}^i$  sont obtenues. Le premier seau reflète en fait le niveau que le système WiMAX est contraint à respecter pour chaque flux conformément au contrat de service ou SLA (Service Level Agreement). Il est à noter que ni la BS ni la MS ne sont contraintes à garantir les délais ( $L_{max}$ ) pour un flux dont le débit dépasse  $R_{min}$ .

Le deuxième seau est utilisé afin de s'assurer que le débit avec lequel le trafic est transmis reste conforme aux limites prédéfinies; i.e. ne dépassant pas  $R_{max}^i$ . Tel que nous le voyons dans la Figure 6.1, le deuxième seau est défini à travers les paramètres suivants:

- un débit moyen appelé "excess information rate" ( $EIR$ ),
- une taille d'excès de rafale  $B_e$
- et un intervalle de temps  $T_e$ .

Nous considérons le même intervalle de mesure que pour le premier seau; i.e.  $T_e = T_c = I_{gr}^i$ . Plus précisément, pour un trafic temps-réel  $i$ ,  $T_e = T_c = L_{max}^i$ .  $B_e$  est configuré de manière à ce que la taille maximale d'une rafale ne dépasse pas  $R_{max}^i \times T_e$ . En d'autres termes,  $B_c + B_e = R_{max}^i \times T_e$  ce qui implique que  $B_e = EIR \times T_e = (R_{max}^i - R_{min}^i) \times T_e$ . Il est à noter que lorsque la capacité  $B_c$  ou  $B_e$  de l'un des seaux est atteinte, le surplus de jetons est supprimé.

Utilisant la configuration décrite ci-dessus, si les seaux sont vides au début de l'intervalle d'allocation, la taille maximale de la rafale ne pourrait être atteinte qu'à la fin de l'intervalle. Plus précisément, si les packets sont générés à un débit  $R_{max}^i$  d'une manière sporadique (toujours conforme au contrat), ils sont automatiquement retardés même si l'on dispose de suffisamment de ressources radio pour les transmettre et ce parce qu'il n'y a pas encore assez de jetons dans les seaux. Cette configuration permet de lisser le trafic et d'éviter les goulots d'étranglement au prochain saut. Néanmoins, ceci pourrait engendrer un gaspillage des ressources et des délais supplémentaires

inutiles.

Pour plus de flexibilité, et pour une plus grande efficacité dans la gestion de la trame physique, nous choisissons de garder les mêmes intervalles de mesures  $T_c$  et  $T_e$ , et les mêmes tailles de rafales  $B_c$  et  $B_e$ . Toutefois, nous considérons les sceaux pleins au début de l'intervalle. Cette configuration, tout en limitant la sporadicité aux seuils souhaités, permet que celle-ci se produise à n'importe quel instant durant l'intervalle de mesure. Il est à noter que pour les connexions BE, le premier sceau est vide étant donné que  $CIR = R_{min}^i = 0$  et pour UGS, c'est le second sceau qui est vide puisque  $R_{max}^i = R_{min}^i$  et  $EIR = R_{max}^i - R_{min}^i$ . Ainsi, la configuration est assez générique pour supporter tous les types de services.

Ce mécanisme de lissage est combiné à une politique d'ordonnancement à deux étapes dont les détails sont fournis dans ce qui suit.

### Un Algorithme d'ordonnancement à Deux Phases

Le processus d'ordonnancement proposé dans cette thèse consiste en trois ordonnanceurs; deux au niveau de la BS: un pour le lien descendant et un autre pour le lien ascendant et un ordonnanceur au niveau de la MS chargé de redistribuer les ressources allouées par la BS entre les différentes connexions UL. Au début de chaque trame, la BS doit décider de la façon dont la bande passante est répartie entre les flux actifs. Le processus d'ordonnancement que nous proposons agit en deux temps. Dans un premier temps, l'objectif est de satisfaire le SLA en garantissant le débit minimum pour les connexions non-BE et les contraintes de délais pour les connexions temps-réel (UGS, erTPS, et rtPS). La fréquence de ces premières allocations est déterminée par l'intervalle d'allocation du flux considéré:  $I_{gr}^i$ . Mappant ceci au mécanisme de seau à jetons, ceci reviendrait à vider le premier seau des flux dont l'intervalle d'allocation expire dans la trame en cours. En procédant de la sorte, nous évitons d'ordonner toutes les connexions à chaque trame ce qui réduirait l'overhead associé à l'accès d'une MS. Les algorithmes correspondant à l'implémentation de cette première phase au niveau de la BS (en DL et UL) et au niveau de la MS sont donnés respectivement par Algorithm 9, Algorithm 11, et Algorithm 10.

Les paramètres considérés dans ces algorithmes sont les suivants:

- $U = \{u1, u2, \dots, uu\}$  l'ensemble des SFs UGS
- $E = \{e1, e2, \dots, ee\}$  l'ensemble des SFs ertPS
- $R = \{r1, r2, \dots, rr\}$  l'ensemble des SFs rtPS
- $N = \{n1, n2, \dots, nn\}$  l'ensemble des SFs nrtPS
- $B = \{b1, b2, \dots, bb\}$  l'ensemble des SFs BE
- $T_f$  : la durée de la trame
- $Gr_1^i$  : la quantité de bande passante allouée à la connexion  $i$  durant la 1<sup>ère</sup> phase du processus d'ordonnancement.
- $Gr_2^i$  : la quantité de bande passante allouée à la connexion  $i$  durant la 2<sup>ème</sup> phase du processus d'ordonnancement.
- $Gr^i$  : la quantité de bande passante allouée à la connexion  $i$  durant tout l'intervalle d'allocation  $I_{gr}^i$ .
- $R_{min}^i$  : le débit minimum réservé pour la connexion  $i$

- $R_{max}^i$  : le débit maximum supporté pour la connexion  $i$
- $L_{max}^i$  : le délai maximum toléré pour la connexion  $i$
- $I_{gr}^i$  : l'intervalle d'allocation pour la connexion  $i$
- $N_q^i$  : le nombre de paquets séjournant dans la file de la connexion  $i$
- $S_q^i$  : la taille de la file de la connexion  $i$
- $t_{cur}$  : le temps système
- $t_{lgr}^i$  : l'instant auquel la connexion  $i$  a reçu la dernière allocation

---

**Algorithm 1:** BS DL Scheduler: 1st round
 

---

**Return:**  $W$  the sum of connections weights to be used in the 2nd round

```

1 Begin
2    $W \leftarrow 0$ 
3   for ( $i = 0; i < 5; i ++$ ) do
4     for ( $j = 0; j < N_{SF}^i; j ++$ ) do
5        $Gr_1^j \leftarrow 0$ 
6        $w^j \leftarrow 0$ 
7       if ( $t_{cur} - t_{lgr}^j \geq I_{gr}^j$ ) then
8          $tmp\_Gr_1^j \leftarrow \min(S_q^j,$ 
9            $R_{min}^j \times I_{gr}^j - Gr^j)$ 
10         $Gr_1^j \leftarrow \text{ovhd\_avail}(tmp\_Gr_1^j, MCS(j))$ 
11         $BW_r \leftarrow BW_r - Gr_1^j$ 
12         $t_{lgr}^j \leftarrow t_{cur}$ 
13         $w^j \leftarrow \min(S_q^j,$ 
14           $R_{max}^j \times I_{gr}^j - Gr^j) - Gr_1^j$ 
15         $Gr^j \leftarrow 0$ 
16         $W \leftarrow W + w^j$ 
17         $W \leftarrow W + \min(S_q^j, R_{max}^j \times I_{gr}^j - Gr^j)$ 
18  return  $W$ 

```

---

Les connexions participant à cette première phase sont considérées dans un ordre de priorité strict: UGS, ertPS, rtPS, et nrtPS. Seul la quantité de données nécessaire à atteindre le débit minimum (en considérant l'overhead correspondant) est allouée. Il est à noter que, du côté de la BS, puisque les flux pourraient être transmis avec différents MCSs, une conversion de la quantité  $Gr_1^i$  en slots ou symboles OFDM est nécessaire afin d'évaluer la bande passante disponible  $BW_r$ , considérée également en slots de temps dans ce cas là (c.f. ligne 10 de Algorithm 9 et ligne 9 de Algorithm 11). Notons également que dans cette stratégie, nous considérons un ratio DL/UL de 1:1 ce qui représente un des ratios typiques recommandés par le forum WiMAX; contrairement au schéma d'ordonnement présenté dans le chapitre 5 où les limites DL/UL sont ajustées de manière dynamique selon les caractéristiques du trafic.

La seconde phase de l'algorithme d'ordonnement est déclenchée par l'éventuelle existence de

---

**Algorithm 2:** SS Scheduler: 1st round**Return:** W the sum of connections weights to be used in the 2nd round

```

1 Begin
2    $W \leftarrow 0$ 
3   for ( $i = 0; i < 5; i++$ ) do
4     for ( $j = 0; j < N_{SF}^i; j++$ ) do
5        $Gr_1^j \leftarrow 0$ 
6        $w^j \leftarrow 0$ 
7       if ( $t_{cur} - t_{lgr}^j \geq I_{gr}^j$ ) then
8          $tmp\_Gr_1^j \leftarrow \min(S_q^j,$ 
9            $R_{min}^j \times I_{gr}^j - Gr^j)$ 
10         $Gr_1^j \leftarrow \text{ovhd\_avail}(tmp\_Gr_1^j)$ 
11         $t_{lgr}^j \leftarrow t_{cur}$ 
12         $w^j \leftarrow \min(S_q^j,$ 
13           $R_{max}^j \times I_{gr}^j - Gr^j) - Gr_1^j$ 
14         $Gr^j \leftarrow 0$ 
15         $W \leftarrow W + w^j$ 
16      else if ( $(i \in R \text{ or } i \in N)$ 
17         $\text{and } (t_{cur} - t_{lgr}^j + T_f \geq I_{gr}^j))$ ) then
18        if ( $\text{unicast\_BR\_Opp} \geq 1$ ) then
19           $\text{send\_standalone\_BR}$ 
20        else if ( $BWr \geq 6$ ) then
21           $\text{/* bandwidth stealing */}$ 
22           $\text{send\_standalone\_BR}$ 
23        else if ( $N_{SF}^0 \geq 1$ ) then
24           $PM\_bit \leftarrow 1$ 
25         $W \leftarrow W + \min(S_q^j, R_{max}^j \times I_{gr}^j - Gr^j)$ 
26  return W

```

**Algorithm 3:** BS UL Scheduler: 1st round**Return:** W the sum of connections weights to be used in the 2nd round

```

1 Begin
2    $W \leftarrow 0$ 
3   for ( $i = 0; i < 5; i++$ ) do
4     for ( $j = 0; j < N_{SF}^i; j++$ ) do
5        $Gr_1^j \leftarrow 0$ 
6       if ( $t_{cur} - t_{lgr}^j \geq I_{gr}^j$ ) then
7          $tmp\_Gr_1^j \leftarrow \min(Req^j,$ 
8            $R_{min}^j \times I_{gr}^j) - Gr^j$ 
9          $Gr_1^j \leftarrow \text{ovhd\_avail}(tmp\_Gr_1^j)$ 
10         $BW_r \leftarrow BW_r - Gr_1^j$ 
11         $t_{lgr}^j \leftarrow t_{cur}$ 
12         $w^j \leftarrow \min(Req^j,$ 
13           $R_{max}^j \times I_{gr}^j) - Gr^j - Gr_1^j$ 
14         $Gr^j \leftarrow 0$ 
15         $W \leftarrow W + w^j$ 
16         $((i \in R \text{ or } i \in N)$ 
17        else if  $\text{and}(t_{cur} - t_{lgr}^j + T_f \geq I_{gr}^j)$  then
18           $\text{and}((N_{SF}^0 == 0)$ 
19           $\text{or}(N_{SF}^0 > 0 \text{ and } PM == 1))$ 
20           $Unicast\_Poll$ 
21           $W \leftarrow W + \min(Req^j, R_{max}^j \times I_{gr}^j - Gr^j)$ 
22        return W

```

bande passante à la fin de la première phase. L'objectif de cette seconde étape est de partager la bande passante restante entre les différentes connexions. Ce partage est assuré conformément à la stratégie WFQ. Le poids de chaque connexion correspond à la taille de sa file d'attente tout en restant dans les limites fixées par le double sceau qui lui correspond. Lorsque  $Gr_2^i$  est calculé, une quantité équivalente de jetons est retirée du premier puis du second sceau.

---

**Algorithm 4: BS DL Scheduler: 2nd round**


---

```

1 Begin
2    $W \leftarrow 0$ 
3   for ( $i = 0; i < 5; i++$ ) do
4     for ( $j = 0; j < N_{SF}^i; j++$ ) do
5        $tmp\_Gr_2^j \leftarrow \frac{w^j}{W} \times BW_r$ 
6        $Gr_2^j \leftarrow ovhd\_avail(tmp\_Gr_2^j)$ 
7        $BW_r \leftarrow BW_r - Gr_2^j$ 
8        $Gr^j \leftarrow Gr^j + Gr_2^j$ 

```

---



---

**Algorithm 5: SS Scheduler: 2nd round**


---

```

1 Begin
2    $W \leftarrow 0$ 
3   for ( $i = 0; i < 5; i++$ ) do
4     for ( $j = 0; j < N_{SF}^i; j++$ ) do
5        $Gr_2^j \leftarrow 0$ 
6       if ( $w^j > 0$ ) then
7          $tmp\_Gr_2^j \leftarrow \frac{w^j}{W} \times BW_r$ 
8          $Gr_2^j \leftarrow ovhd\_avail(tmp\_Gr_2^j)$ 
9          $BW_r \leftarrow BW_r - Gr_2^j$ 
10         $Gr^j \leftarrow Gr^j + Gr_2^j$ 
11       if ( $Gr_2^j > 0$  and  $S_q^j > 0$ ) then
12         if ( $BW_r > 2$ ) then
13            $Piggyback\_BR$ 
14         else if ( $Contention\_BR\_Opp$ ) then
15            $send\_standalone\_BR$ 

```

---

Dans cette seconde phase, les connexions BE se voient accorder, proportionnellement, autant de chances que d'autres types de flux pour concourir pour une partie de la bande passante ce qui éviterait des problèmes de famine. Les détails des algorithmes proposés sont fournis dans Algorithm 12 (coté BS en DL) et Algorithm 13 (coté MS), respectivement. La Figure 4 illustre les trois cas de figure possibles des sceaux à jetons à la fin d'un intervalle d'allocation pour une connexion donnée  $i$ , après avoir effectué les deux étapes d'ordonnancement. Il est à noter que durant tout l'intervalle, aucun jeton n'est rajouté.

---

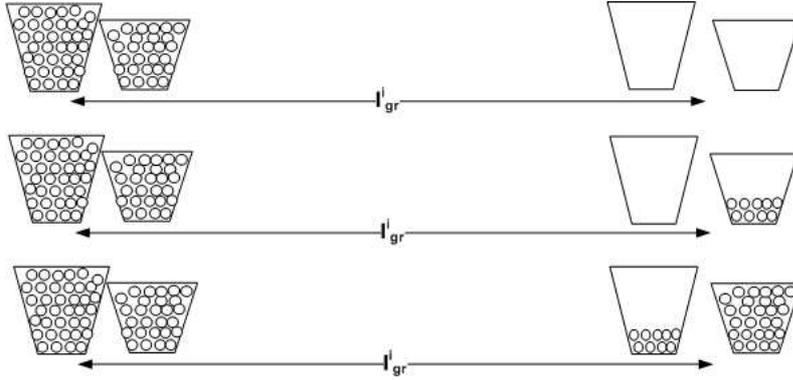


Figure 4: A Dual-Bucket Shaping Mechanism

- Dans le premier cas, les deux seaux sont vides ce qui implique que la connexion a atteint son débit maximum  $R_{max}^i$ .
- Si seul le premier seau est vide, cela veut dire que la connexion a été ordonnancé avec un débit  $R^i$ ;  $R_{min}^i \leq R^i < R_{max}^i$ . Ce qui veut dire que la connexion a réussi à atteindre au moins le débit minimum garanti et à ne pas dépasser le délai maximum toléré.
- Le troisième cas illustrée par la figure 4 correspond au cas où le premier seau n'est pas complètement vidé i.e.  $R^i < R_{min}^i$ . En d'autres termes, la bande passant n'était pas suffisamment large pour couvrir les besoins des connexions participant à la première phase de l'algorithme.

Dans les deux premiers cas, les deux seaux associés au flux sont remplis de nouveau et un nouvel intervalle d'allocation commence. Dans le dernier cas, par contre, les mêmes seaux sont maintenus. De plus, pour atteindre  $R_{min}^i$ , la connexion a besoin de plus de bande passant que ce qui est reflété par le contenu en jetons de son premier seau. De ce fait, au début de la trame d'après,  $T_f \times R_{min}^i$  jetons du deuxième seau sont marqués indiquant que le seuil pour la première étape ne correspond pas uniquement au contenu du premier seau mais également aux jetons marqués du deuxième seau. La connexion participe à la première phase autant de fois que nécessaire, durant les prochaines trames, jusqu'à ce que tous les jetons du premier seau ainsi que ceux marqués du deuxième seau soient utilisés. C'est d'ailleurs seulement à cet instant que les seaux sont de nouveau remplis de jetons. Ce dernier cas entraînerait un délai supplémentaire pour le flux considéré. Néanmoins, en décalant l'intervalle d'allocation nous diminuons les chances que ce cas de figure se produise encore une fois (deux ou plusieurs rafales coincident), surtout si cette sporadicité survient d'une manière périodique.

### Analyse de Performance

Afin d'évaluer les performances de la stratégie mCoSS, nous avons implémenté l'ensemble des algorithmes associés sous Qualnet 4.5 [31], qui est la version commerciale de GloMoSim. mCoSS a été comparée aux disciplines SP et à une variante du WFQ.

Le Tableau 2 dresse les paramètres de simulation considérés dans notre évaluation de performances.

Dans les scénarios qui suivent, nous considérons un stream audio de 30 mns configuré comme une connexion UL du type rtPS. La taille de la trame audio est fixée à 1600 octets et le nombre de

Fréquence du canal	3.5 GHz
Band passante	10 MHz
Taille FFT	2048
Gain du préfixe cyclique	8
Modèle de propagation	Two-ray
Puissance d'émission de l'antenne de la BS	33 dBm (= 2 W)
Hauteur de l'antenne de la BS	32 m
Gain de l'antenne de la BS	15 dBi
Puissance d'émission de l'antenne de la MS	23 dBm (= 200 mW)
Hauteur de l'antenne de la MS	1.5 m
Gain de l'antenne de la MS	-1 dBi
Type d'antenne	omnidirectionnel
Durée de la trame	10 ms
Durée de la portion DL	5 ms

Table 2: Simulation settings

trames par secondes suit une distribution uniforme entre 10 et 25 (trames/s). Les paramètres de QoS considérés pour ce flux audio sont les suivants:  $R_{min}^i = 128$  kbps,  $R_{max}^i = 320$  kbps et  $I_{gr}^i = 100$  ms.

### Scénario 1

A travers ce scénario, nous nous proposons d'évaluer le pouvoir de shaping de notre stratégie mCoSS. Pour cela, nous plaçons deux MSs à distance égale d'une station de base et nous configurons le stream audio comme mentionné précédemment:  $R_{min}^i = 128$  kbps,  $R_{max}^i = 320$  kbps et  $I_{gr}^i = 100$  ms. Tandis que MS1 respecte ces limites, MS2 tente de transmettre à un débit beaucoup plus élevé variant de 640 kbps à 1.28 Mbps. Plus de 30 expériences ont été tournées afin de valider la capacité de mCoSS à lisser un trafic gourmand en bande passante et à comparer son comportement à celui d'une variante du WFQ et au SP implémentés sous Qualnet.

La Figure 5 représente les débits d'émission et de réception des deux flux: celui conforme au contrat (well-behaving) et celui qui est gourmand (misbehaving) pour les trois algorithmes: mCoSS, WFQ et SP. Le "Tx rate" représente le débit avec lequel l'application est générée au niveau de la MS tandis que le "Rx Rate" est le débit de réception à la BS.

Nous pouvons observer à partir de la Figure 5 que pour le trafic conforme envoyé par MS1 les trois stratégies ont des performances quasi identiques en terme de débit. Pour le trafic gourmand en bande passante, SP et WFQ laisse le trafic atteindre plus de 800 kbps alors que mCoSS oblige ce trafic à rester dans les limites fixées par le contrat de QoS; le débit à la réception ne dépasse pas en effet les 315 kbps.

Les tableaux 3 et 3 reportent les valeurs obtenues pour les délais de bout en bout et la gigue pour les deux trafics observés. Comme conséquence de la politique de shaping adoptée par notre stratégie, le trafic non conforme au contrat généré par MS2 est pénalisé par mCoSS (en comparaison à WFQ et SP) en terme de délai puisque les packets dépassant  $R_{max}^i$  sont retardés et éventuellement supprimés si leur nombre dépasse la capacité des buffers. D'autre part, les délais du trafic conforme sont réduits de moitié (en comparaison à WFQ et SP). Avec WFQ et SP, les deux trafics obtiennent les mêmes délais de bout en bout; le trafic gourmand bénéficie même d'une gigue moyenne plus courte que celle du trafic conforme. A partir des résultats obtenus, nous pouvons constater que mCoSS est capable de forcer un trafic à rester dans la limite des seuils autorisés et d'isoler les

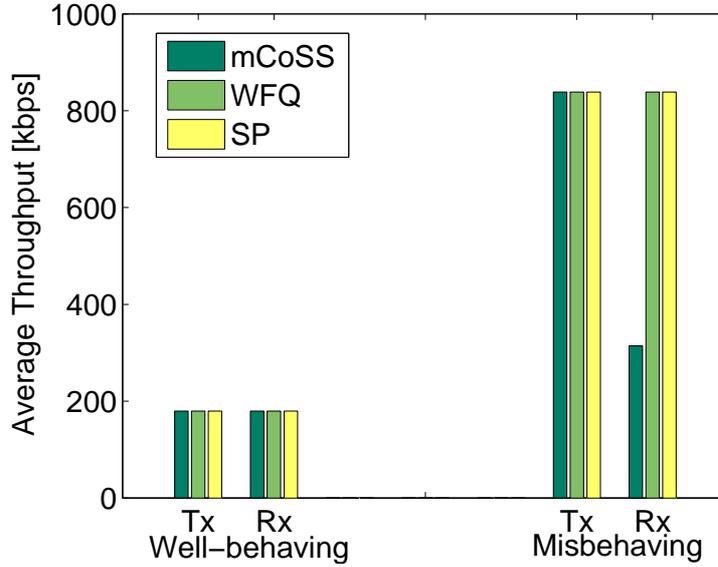


Figure 5: mCoSS Shaping Capability

	MS1 Well-behaving	MS2 Misbehaving
mCoSS	0.255	13.6
WFQ	0.57	0.53
SP	0.57	0.53

Table 3: mCoSS Shaping Capability: E2E Delay (sec)

trafics gourmands. L'absence d'un mécanisme de shaping dans WFQ et SP a affecté les performances du premier trafic et les conséquences auraient pu être plus significatives si le second trafic avait tenté de saturer toute la bande passante.

## Scénario 2

Dans ce second scénario, nous considérons les mêmes MSs avec chacune trois streams audio ayant la même configuration. A travers ce scénario, nous visons à évaluer, dans des conditions équivalentes de canal et de trafic, les performances de notre stratégie d'ordonnancement en terme d'équité inter-SSs et inter-SFs et de comparer le degré de satisfaction en QoS des six connexions en utilisant les trois stratégies d'ordonnancement. La Figure 6(a) montre le débit moyen obtenu pour le 1er, 2nd et 3ème stream audio (A1, A2 et A3, respectivement) de MS1 et MS2. Coté débit moyen, les trois stratégies offrent le même niveau de performance. Le délai et la gigue de bout en

	MS1 Well-behaving	MS2 Misbehaving
mCoSS	22	80
WFQ	69	27.7
SP	69	27.7

Table 4: mCoSS Shaping Capability: Jitter (ms)

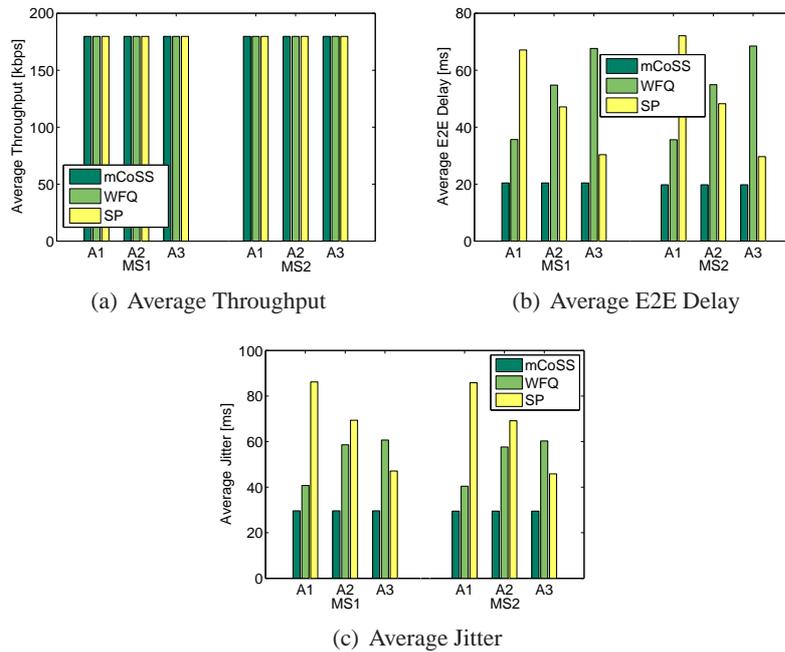


Figure 6: 2 MSs with 3 Audio streams each

bout quant à eux connaissent des comportements plus variables d'un algorithme à un autre comme nous pouvons le constater à partir des Figures 6(b) and 6(c). Avec le WFQ, le délai E2E varie de 35 à 67 ms d'un SF à un autre. Le même comportement est observé pour SP pour lequel le délai moyen de bout en bout varie de 30 à 72 ms. mCoSS d'un autre coté offre des valeurs bien plus basses et beaucoup plus stables pour les six flux aussi bien pour le délai (environ 20 ms) que pour la gigue (moins de 30 ms). Comparé au WFQ et au SP, mCoSS offre les meilleurs et surtout les plus stables résultats ce qui entraîne une meilleure équité inter-SFs et inter-MSs.

### Scénario 3

A travers ce dernier scénario, nous tentons de valider la capacité de mCoSS à adapter la bande passante allouée aux conditions du canal de la MS; une fonctionnalité qui est déjà supportée par le module WiMAX dans Qualnet pour les algorithmes WFQ et SP.

Etant donné cet objectif, nous considérons 3 MSs placées à des positions plus ou moins éloignées de la BS: à 1km, 2 km et 3km. Ces trois distances correspondent en fait à trois niveaux de SNR correspondant à UIUC 1 (QPSK 1/2), UIUC 4 (16-QAM 3/4) and UIUC 7 (64-QAM 3/4). Nous configurons deux streams audio à chaque MS avec les mêmes paramètres spécifiés précédemment. Tel que le montre la Figure 7(a), d'ailleurs comme pour le scénario précédent, les trois algorithmes ont des performances quasi équivalentes pour ce qui est du débit moyen. Cependant, la différence du délai moyen de bout en bout (illustré dans la Figure 7(b)) entre Audio 1 et Audio 2, en utilisant la stratégie SP est plus visible que dans le cas du scénario précédent. En effet il varie par exemple pour MS3 de 35 ms à plus de 100 ms excédant ainsi le délai maximum toléré. Ce même comportement est observé pour la gigue moyenne dans la Figure 7(c). Pour mCoSS par contre, l'utilisation de différents schémas de modulation et de codage n'a pratiquement eu aucun effet sur les performances de l'algorithme. L'équité et la stabilité des résultats observées dans le scénario précédent sont confirmées à travers ce scénario.

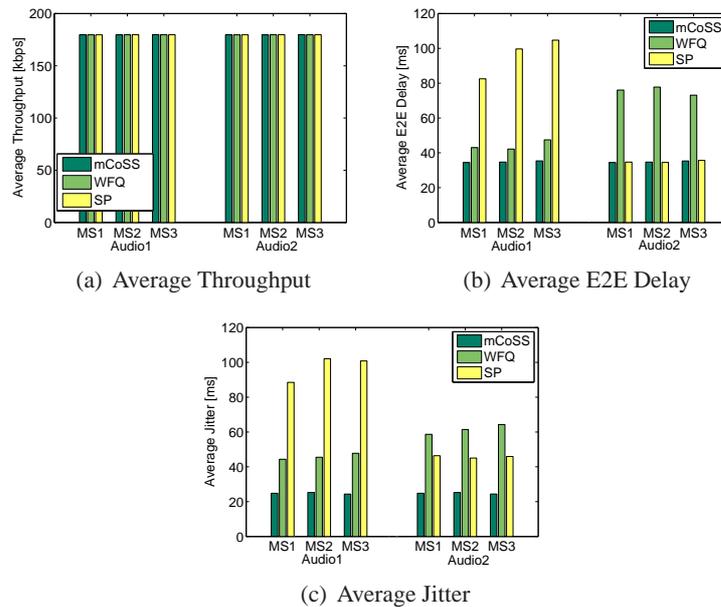


Figure 7: 3 MSs with 2 Audio streams each

### Conclusion et perspectives

La majorité des stratégies d'ordonnancement hiérarchiques (telles que [13, 10, 9]) proposées dans la littérature et décrites dans le Chapitre 4 propose une discipline de gestion de file d'attente spécifique à chaque service d'ordonnancement (e.g. EDF pour rtPS, WFQ pour nrtPS et RR pour BE). Ceci entraîne une augmentation significative de la complexité de la stratégie adoptée. Contrairement à ces approches, la stratégie mCoSS que nous proposons dans cette thèse est conçue de manière à pouvoir s'appliquer à tous les types de services d'ordonnancement visés par la technologie WiMAX.

Basée sur une variante du mécanisme du double seau à jetons, mCoSS allie la généralité de l'approche à la spécificité de la configuration dans la mesure où ce mécanisme de shaping est configuré par flux.

La politique de demande et d'allocation de la bande passante adoptée par mCoSS est conçue de manière à assurer un compromis entre la précision de la perception des besoins en bande passante et la diminution de l'overhead associé à un polling unicast plus fréquent. En effet la stratégie alterne "bandwidth stealing", piggybacking, polling unicast ou broadcast et usage du PM bit en fonction du service d'ordonnancement considéré et des ressources disponibles.

Les résultats préliminaires d'évaluation de mCoSS reportés dans ce manuscrit valident et confirment l'équité de la stratégie, sa capacité à isoler et lisser les trafics gourmands en bande passante et à supporter la technique très intéressante d'AMC qu'offre le WiMAX. Plus de simulations sont toutefois nécessaires pour vérifier et valider d'autres aspects de cette stratégie.

### Sujets Relatifs à la Gestion de la Mobilité dans les Réseaux WiMAX

Les principales conclusions dégagées au niveau de cette annexe pourraient être résumées comme suit:

- Le standard IEEE 802.16e propose trois modes de handover. Le hard handover consiste en

plus d'étapes et pourraient entraîner des délais importants. Néanmoins, les deux modes de soft handover: FBSS et MDHO ne peuvent pas constituer une alternative fiable au schéma de handover obligatoire (hard HO) pour plusieurs raisons. D'une part parce qu'il existe plusieurs restrictions sur les BSs opérant en modes FBSS/MDHO étant donné qu'elles doivent être synchronisées entre elles en temps et en fréquence et doivent avoir des structures de trames synchrones. D'autre part, dans les modes MDHO et FBSS, les BSs d'un même Diversity Set appartiendraient a priori à un même sous réseau tandis que le handover pourraient survenir entre BSs appartenant à des sous réseaux différents.

- Les handovers rapides de Mobile IPv6, comme tous les mécanismes de gestion des handovers cross-layer en général, sont basés sur la collaboration de différentes couches afin d'améliorer la manière dont la mobilité est gérée. Cette idée d'intégrer les informations provenant de plusieurs couches aide à augmenter les performances de gestion du HO. Néanmoins, ces solutions requièrent souvent d'importantes modifications chose qui complique et freine leur déploiement.
  - Les réseaux émergents seraient a priori hétérogènes; donc, la convergence vers un mécanisme de handover unifié est désormais une nécessité. De ce fait, le mécanisme Media Independent Handover (MIH) offre une alternative intéressante dans la mesure où il offre une solution généralisée et surtout standardisée pour différentes technologies d'accès. Toutefois, le succès du MIH dépend énormément de la bonne volonté des vendeurs à l'intégrer dans leur produits futurs.
  - Le roaming constitue un concept clef pour l'extension de la couverture d'un réseau d'opérateurs. A travers le roaming, un mobile pourrait accéder automatiquement aux services d'un autre opérateur lorsqu'il se trouve en dehors de la couverture du réseau de son opérateur habituel (home network provider). Le roaming propose de ce fait un modèle économique plus générique et plus extensible pour les réseaux WiMAX. D'ailleurs, afin d'offrir un processus indépendant et plus étendu, le forum WiMAX définit une interface de roaming. L'objectif de cette interface est de standardiser le format ainsi que les moyens d'échange entre les différentes entités impliquées dans le processus de roaming.
-



# Acknowledgements

I would like first to thank Dr. Fethi Filali, my supervisor, for offering me the opportunity to pursue my doctoral studies at EURECOM, for his insightful guidance, and for being available whenever I needed his support.

My deepest gratitude goes to Daniel Câmara for the collaboration we had: many thanks Daniel for all the time you have dedicated to our discussions.

My special thanks go to my great friends Randa, Zuleita, Saoucene, Sara, and Giuliana for cheering me up whenever I was down. Not forgetting the nice environment and support provided by my office mates and friends Erhan, Kostas, Bassem, Aymen, Umer, Antony, and Mustapha.

I would like to extend my thanks to our Department Head Prof. Christian Bonnet, to the secretaries, and to the IT Department staff for all the efforts they are putting to offer a productive research environment at EURECOM.

I would never thank enough my father Hedi and my mother Samira for their love, trust, and support.

Last but not least, I am grateful to my little Cyrine and to my dear husband Sami for putting up with my lack of availability. I hope finishing my doctoral studies would help me free up more time for them.

---



---

# Contents

<b>List of Figures</b>	<b>33</b>
<b>Acronyms</b>	<b>35</b>
<b>Introduction</b>	<b>39</b>
<b>1 An overview of WiMAX</b>	<b>43</b>
1.1 IEEE 802.16/WiMAX standardization . . . . .	43
1.2 IEEE 802.16 frequency spectrum and PHY interfaces . . . . .	44
1.3 An overview of WiMAX PHY . . . . .	44
1.3.1 Link adaptation, modulation, and coding . . . . .	45
1.3.2 WirelessMAN-OFDM . . . . .	46
1.3.3 WirelessMAN-OFDMA . . . . .	47
1.3.3.1 Subchannelization . . . . .	47
1.3.3.2 Band AMC Permutation mode . . . . .	47
1.3.3.3 OFDMA frame structuring . . . . .	47
1.4 An overview of WiMAX MAC . . . . .	48
1.4.1 A connection-oriented MAC . . . . .	48
1.4.2 IEEE 802.16 protocol stack . . . . .	49
1.4.3 MAC PDUs formats, construction and transmission . . . . .	50
1.5 Conclusion . . . . .	51
<b>2 Performance Analysis of OFDM-based WiMAX Networks</b>	<b>53</b>
2.1 Analytical framework . . . . .	53
2.2 Performance evaluation . . . . .	58
2.2.1 Effect of the frame duration and the MCS . . . . .	58
2.2.2 Effect of the channel bandwidth . . . . .	62
2.2.3 Impact of fragmentation and packing . . . . .	62
2.3 Conclusion . . . . .	64
<b>3 QoS Support in WiMAX Networks</b>	<b>67</b>
3.1 QoS support in WiMAX networks . . . . .	67
3.1.1 Service flows management and QoS requirements . . . . .	67
3.1.2 Scheduling service types . . . . .	68
3.1.3 Bandwidth allocation and request mechanisms . . . . .	70
3.2 A QoS architecture for WiMAX networks: the big picture . . . . .	71
3.3 Scheduling and CAC in WiMAX: design challenges . . . . .	75
3.4 Conclusion . . . . .	76

---

---

<b>4</b>	<b>Scheduling and CAC in WiMAX Networks: a Survey and Taxonomy</b>	<b>77</b>
4.1	Scheduling . . . . .	77
4.1.1	Packet queuing-derived strategies . . . . .	77
4.1.1.1	One-layer scheduling structures . . . . .	77
4.1.1.2	Hierarchical scheduling structures . . . . .	79
4.1.2	Optimization-based strategies . . . . .	86
4.1.3	Cross-layer strategies . . . . .	87
4.2	CAC . . . . .	90
4.2.1	CAC schemes with degradation strategy . . . . .	90
4.2.1.1	Service degradation . . . . .	91
4.2.1.2	Bandwidth borrowing . . . . .	91
4.2.1.3	Bandwidth stealing . . . . .	92
4.2.2	CAC schemes without degradation strategy . . . . .	93
4.2.3	Other CAC schemes . . . . .	93
4.2.3.1	AMC-induced CAC: . . . . .	94
4.2.3.2	CAC for real-time video applications: . . . . .	94
4.3	Conclusion . . . . .	94
<b>5</b>	<b>Adaptive Scheduling with Max-Min Fairness Admission Control</b>	<b>97</b>
5.1	Uplink and downlink scheduling . . . . .	97
5.1.1	Hierarchical scheduling structure . . . . .	97
5.1.2	The BS scheduling algorithm . . . . .	98
5.1.2.1	Step 1: Initialize the available time . . . . .	99
5.1.2.2	Step 2: Plan the first burst . . . . .	99
5.1.2.3	Step 3: Proceed in accordance with the scheduling structure . . . . .	100
5.1.2.4	Step 4: Share bandwidth and plan transmissions . . . . .	103
5.1.3	Admission control policy . . . . .	103
5.2	Performance analysis . . . . .	105
5.2.1	Scenario 1: Single SS scenario . . . . .	105
5.2.2	Scenario 4: Multiple SSs scenario . . . . .	107
5.3	Conclusion . . . . .	108
<b>6</b>	<b>mCoSS: a multi-Constraints Scheduling Strategy for WiMAX Networks</b>	<b>111</b>
6.1	A modified dual-bucket shaping mechanism . . . . .	112
6.2	A two-rounds scheduling algorithm . . . . .	113
6.2.1	Bandwidth request and grant strategy . . . . .	119
6.3	Performance Analysis . . . . .	119
6.3.1	A WiMAX simulation model under QualNet . . . . .	119
6.3.2	Performance evaluation . . . . .	120
6.3.2.1	Scenario 1: mCoSS shaping capability . . . . .	120
6.3.2.2	Scenario 2: fairness and QoS degree of satisfaction . . . . .	122
6.3.2.3	Scenario 3: AMC support . . . . .	123
6.4	Conclusion . . . . .	123
<b>7</b>	<b>Mobile WiMAX: a V2I Communications Medium</b>	<b>125</b>
7.1	ITS applications and architectures . . . . .	125
7.2	IEEE 802.11p vs. IEEE 802.16e . . . . .	127
7.3	Performance evaluation . . . . .	128

---

---

7.3.1	Simulation environment and settings . . . . .	128
7.3.2	Performance analysis . . . . .	131
7.3.2.1	Scenario 1: Impact of the source data rate on the performance . . . . .	131
7.3.2.2	Scenario 2: Impact of the vehicle speed on the performance . . . . .	132
7.4	Conclusion . . . . .	133
<b>Conclusion</b>		<b>134</b>
<b>A Topics Related to Mobility Management in WiMAX Networks</b>		<b>139</b>
A.1	Mobile WiMAX architecture . . . . .	139
A.2	Horizontal handover in 802.16e . . . . .	141
A.2.1	Network topology acquisition . . . . .	142
A.2.2	Handover process . . . . .	145
A.2.2.1	Cell reselection . . . . .	145
A.2.2.2	HO decision and initiation . . . . .	145
A.2.2.3	Synchronization to target BS downlink . . . . .	145
A.2.2.4	Ranging and network re-entry . . . . .	145
A.2.2.5	Termination of MS context . . . . .	146
A.2.3	Fast BS switching (FBSS) and macro diversity handover (MDHO) . . . . .	146
A.2.3.1	Macro diversity handover (MDHO) . . . . .	146
A.2.3.2	Fast BS switching (FBSS) . . . . .	147
A.3	Optimized 802.16e handover schemes . . . . .	147
A.3.1	L2 handover schemes . . . . .	147
A.3.2	L2-L3 cross-layer handover schemes . . . . .	149
A.3.3	Mobile IPv6 fast handovers over IEEE 802.16e networks . . . . .	150
A.3.3.1	Predictive mode . . . . .	150
A.3.3.2	Reactive mode . . . . .	151
A.4	Vertical handover . . . . .	153
A.4.1	Vertical handover mechanisms involving 802.16e networks . . . . .	154
A.4.2	IEEE 802.21, media-independent handover services . . . . .	155
A.4.2.1	General architecture . . . . .	155
A.4.2.2	MIHF services . . . . .	156
A.5	Roaming . . . . .	158
A.6	Conclusion . . . . .	159
<b>Bibliography</b>		<b>166</b>

---



# List of Figures

1	QoS architecture Design . . . . .	11
2	Classification of the scheduling strategies of IEEE 802.16 PMP mode . . . . .	12
3	A Dual-Bucket Shaping Mechanism . . . . .	14
4	A Dual-Bucket Shaping Mechanism . . . . .	20
5	mCoSS Shaping Capability . . . . .	22
6	2 MSs with 3 Audio streams each . . . . .	23
7	3 MSs with 2 Audio streams each . . . . .	24
1.1	OFDM Frame Structure with TDD . . . . .	47
1.2	Example of OFDMA frame in TDD . . . . .	48
1.3	Slot structure in Band AMC Permutation . . . . .	49
1.4	IEEE Std 802.16 Data Plane Protocol Reference Model . . . . .	50
1.5	IEEE 802.16 MAC Headers Formats . . . . .	51
2.1	OFDM Symbol structure . . . . .	54
2.2	OFDM Frame Structure with TDD . . . . .	55
2.3	Effect of frame duration and modulation and coding scheme on IP throughput . . . . .	60
2.4	Effect of the frame duration and the MCS on bandwidth utilization . . . . .	61
2.5	Effect of the channel bandwidth on MAC efficiency . . . . .	63
2.6	MAC efficiency using 64-QAM 3/4 . . . . .	64
2.7	Effect of packing and fragmentation . . . . .	65
3.1	Dynamic Service Addition . . . . .	69
3.2	QoS architecture Design . . . . .	72
4.1	Classification of the scheduling strategies of IEEE 802.16 PMP mode . . . . .	78
4.2	Hierarchical structure for bandwidth allocation [13, 14] . . . . .	80
4.3	3 schedulers proposal[8] . . . . .	81
4.4	Multimedia supported uplink scheduler [10] . . . . .	81
4.5	Scheduler model for WiMAX[11] . . . . .	82
4.6	Hierarchical structure of bandwidth allocation for WiMAX PMP mode [7] . . . . .	83
4.7	Operation flowchart of 2TSA [6] . . . . .	86
5.1	Min-Max CAC Policy . . . . .	106
5.2	Single SS scenario . . . . .	107
5.3	Multiple SSs Scenario . . . . .	109
6.1	A Dual-Bucket Shaping Mechanism . . . . .	113
6.2	A Dual-Bucket Shaping Mechanism . . . . .	118
6.3	mCoSS Shaping Capability . . . . .	121

---

6.4	2 MSs with 3 Audio streams each . . . . .	122
6.5	3 MSs with 2 Audio streams each . . . . .	123
7.1	ITS station reference architectures. . . . .	127
7.2	European channel allocation [32]. . . . .	128
7.3	Coverage evaluation scenarios. . . . .	130
7.4	Scenarios network deployments. . . . .	131
7.5	Impact of the source data rate on the average performance . . . . .	132
7.6	Impact of the vehicle speed on the average performance . . . . .	133
A.1	Network Reference Model . . . . .	140
A.2	ASN interoperability Profiles [33, 34] . . . . .	140
A.3	Example of neighbor BS advertisement and scanning (without association) by MS request [2] . . . . .	142
A.4	Example of neighbor BS advertisement and scanning (with non-coordinated association) by MS request [2] . . . . .	143
A.5	Example of macro diversity HO (Diversity Set Update: Add) [2] . . . . .	148
A.6	Example of a handover between two different subnets . . . . .	150
A.7	Predictive fast handover in 802.16e [35] . . . . .	152
A.8	Reactive fast handover in 802.16e [35] . . . . .	153
A.9	MIH Reference Model and Services [36] . . . . .	155

---

---

# Acronyms

Here are the main acronyms used in this document. The meaning of an acronym is usually indicated once, when it first occurs in the text.

2TSA	two-tier scheduling algorithm
AF	assured forwarding
AMC	adaptive modulation and coding
APA	adaptive power allocation
ARQ	automatic repeat request
AWGN	additive white gaussian noise
BE	best effort
BPSK	binary phase shift keying
BR	bandwidth request
BS	base station
BWA	broadband wireless access
BWN	broadband wireless network
CAC	connection admission control
CBR	constant bit rate
CIR	committed information rate
CDMA	code division multiple access
CDC	combined distributed and centralized
CID	connection identifier
CL	controlled load
CPE	consumer premises equipment
CPS	common part sublayer
CQI	channel quality information
CRC	cyclic redundancy check
CRRM	common radio resource manager
CS	convergence sublayer
CTMC	continuous time markov chain
DCD	downlink channel descriptor
DFPQ	deficit fair priority queuing
DHCP	dynamic host configuration protocol
DiffServ	differentiated services
DIUC	downlink interval usage code
DL	downlink
DLFP	downlink frame prefix
DRR	deficit round robin
DSA	dynamic service addition
DSC	dynamic service change

---

---

DSD	dynamic service deletion
EDD	earliest due date
EDF	earliest deadline first
EF	expedited forwarding
ertPS	extended real-time polling service
ETSI	european telecommunications standards institute
FDD	frequency division duplex or duplexing
FDMA	frequency division multiple access
FEC	forward error correction
FIFO	first in first out
FQ	fair queuing
FTP	file transfer protocol
GM	grant management
GOP	group of pictures
GPC	grant per connection
GPSS	grant per subscriber station
GS	guaranteed service
HO	handover
HTTP	hypertext transfer protocol
IE	information element
IP	Internet protocol
IEEE	institute of electrical and electronics engineers
IMT	international mobile telecommunications
IntServ	integrated services
ISP	Internet service provider
IUC	interval usage code
L2	layer 2
L3	layer 3
LOS	line-of-sight
LR	latency-rate
LST	latest starting time
LTE	Long Term Evolution
MAC	media access control
MCS	modulation and coding scheme
MIMO	multi-input multi-output
MMFS	max-min fair sharing
MPEG	moving picture experts group
NLOS	non-line-of-sight
nrtPS	non-real-time polling service
OFDM	orthogonal frequency division multiplexing
OFDMA	orthogonal frequency division multiple access
OSI	open systems interconnection
PDRR	pre-scale dynamic resource reservation
PDU	protocol data unit
PF	proportional fair
PHS	payload header suppression
PHY	physical layer
PKM	key management protocol

---

---

PMP	point-to-multipoint
PQLW	priority-based queue length weighted
QAM	quadrature amplitude modulation
QoS	quality of service
QPSK	quadrature phase-shift keying
RED	random early detection
RF	radio frequency
RR	round-robin
RRM	radio resources management
RS-CC	Reed–Solomon-convolutional code
RTG	receive/transmit transition gap
rtPS	real-time polling service
SAP	service access point
SAQoS	service adaptive quality of service
SC	single-carrier
SCFQ	self-clocked fair queuing
SD	silence detector
SF	service flow
SFID	service flow identifier
SINR	signal-to-interference-noise ratio
SLA	service level agreement
SMTP	simple mail transfer protocol
SNMP	simple network management protocol
SNR	signal-to-noise ratio
SPLF	shortest packet length first
SP-order	shortest path order
SS	subscriber station
TAC	threshold-based admission control
TCP	transmission control protocol
TDD	time division duplex or duplexing
TDM	time division multiplexing
TDMA	time division multiple access
TGd	task group d
TGe	task group e
TLV	type-length-value
TTG	transmit/receive transition gap
UCD	uplink channel descriptor
UDP	user datagram protocol
UGS	unsolicited grant service
UIUC	uplink interval usage code
UL	uplink
VAD	voice activity detection
VoIP	voice over IP (Internet protocol)
WFQ	weighted fair queuing
WG	working group
WiMAX	worldwide interoperability for microwave access
WiMesh	wireless mesh
WirelessMAN	wireless metropolitan area networks

---

WMN	wireless mesh networks
WRR	weighted round-robin

---

# Introduction

## Motivation

Over the past two decades, our daily lives have been reshaped by the fast development in the telecommunications environment. Broadband Internet and wireless ubiquity have become more than ever real needs in our modern lifestyle. Driven by this growing demand for high-speed broadband wireless services, Worldwide Interoperability for Microwave Access (WiMAX) technology, addressed in this thesis, has been developed. In addition to its wireless and broadband capability, WiMAX technology is IP-based and mobile. The support of these features makes Mobile WiMAX the leading technology that overcomes the high data fees of 3G technologies data services and the limited mobility of WiFi. Moreover, Mobile WiMAX is a reality and is being deployed in the United States, Japan, Korea, Europe, Australia, and around the globe. It is actually the only mobile broadband technology currently in use. More importantly, there are ongoing discussions about the possible selection of this technology as an International Mobile Telecommunications (IMT)-advanced standard.

WiMAX technology is based on IEEE 802.16 standards and amendments specifying the MAC and PHY layers for fixed, nomadic, portable, and mobile access. The technology offers a set of key features: (1) the use of orthogonal frequency division multiplex (OFDM), (2) time and frequency duplex (TDD and FDD), (3) support of adaptive modulation and coding (AMC) and (4) advanced antenna techniques such as multiple input, multiple output (MIMO) antenna, (5) robust security and (6) Quality-of-Service (QoS) support. In this thesis, we are mainly interested in the latter capability. In fact, WiMAX technology is designed to support heterogeneous classes of services including data, voice and video. However, the IEEE 802.16 standard leaves unstandardized the resource management and scheduling mechanisms which are crucial components to guarantee QoS performance.

In this thesis<sup>2</sup>, we evaluate the performance of WiMAX networks in both fixed and highly mobile environments. More specifically, we investigate the potential and limitations of using Mobile WiMAX as a Vehicle-to-Infrastructure communication medium. Moreover, we tackle in this thesis most of the resource management and scheduling issues that have been left open with the objective of defining an architecture that fulfills the QoS expectations of the five categories of applications addressed by the IEEE 802.16 standard. The remainder of the thesis is organized as described in next section.

---

<sup>2</sup>This work was supported by WiNEM (WiMAX Network Engineering and Multihoming) project under the grant No. 2006 TCOM005 05 and by EURECOM industrial members : BMW Group, Cisco, Monaco Telecom, Orange, SAP, SFR, Sharp, STEricsson, Swisscom, Symantec, Thales.

---

## Contributions and Outline

### Chapter 1: An overview of WiMAX

The objective of this chapter is to provide a broad view of WiMAX technology. Therefore, we first go through the standardization process of the IEEE 802.16 family of standards. Then, we describe the different physical interfaces targeted by the IEEE Std 802.16 as well as the frequency bands for which they have been specified. An overview of the PHY layer is provided with a particular insight into the adaptive modulation and coding capability supported by WiMAX. As for the MAC layer, only the core functionality, necessary to the understanding of the performance study carried out in Chapter 2, is described. All the features related to QoS support at the MAC level are further discussed in Chapter 3. Indeed, because there are so many concepts to be introduced in this context, we have preferred to dedicate a whole chapter to this purpose.

### Chapter 2: Performance Analysis of OFDM-based WiMAX Networks

In this chapter, we evaluate the performance bounds of WiMAX systems under different physical and MAC parameters settings. The saturation throughput that can be reached in 802.16 networks is investigated through several scenarios in which we vary for instance the frame duration, the channel bandwidth, and the modulation and coding scheme (MCS) in use. An analytical framework was developed based on technical properties and system profiles specified by the IEEE 802.16 standard for systems using the WirelessMAN-OFDM air interface.

Parts of this chapter were published in:

- Ikbal Chammakhi Msadaa and Fethi Filali. On the Performance Bounds of OFDM-based 802.16 Broadband Wireless Networks. In WCNC 2008, IEEE Wireless Communications and Networking Conference, Apr. 2008.

### Chapter 3: QoS Support in WiMAX Networks

The IEEE 802.16 standard defines a connection-oriented MAC protocol that is designed to accommodate a variety of applications with different QoS requirements. Nevertheless, several issues mainly related to resource allocation, have been left open. The main objective of this chapter is to provide a better understanding of the supported and missing features to ensure QoS support in WiMAX networks. Therefore, we first describe the main elements specified by the IEEE 802.16 standard to provide QoS for heterogeneous classes of traffic. Then, we propose a generic QoS framework which incorporates what we consider as key components to handle QoS in WiMAX systems. The last section of this chapter is dedicated to scheduling and admission control issues. More specifically, we highlight, in that section, the main challenges faced when designing a scheduling and/or connection admission control (CAC) solution for WiMAX networks.

Parts of this chapter were published in:

- Ikbal Chammakhi Msadaa, Fethi Filali, and Farouk Kamoun. An 802.16 Model for NS2 Simulator with an Integrated QoS architecture. In SIMUTools' 08, 1st International Conference on Simulation Tools and Techniques for Communications, Networks and Systems, Mar. 2008.
-

#### **Chapter 4: Scheduling and CAC in WiMAX Networks: a Survey and Taxonomy**

A large body of literature has been concerned with scheduling and admission control issues in WiMAX networks. In this chapter, we survey, classify, and compare different scheduling and CAC mechanisms proposed in this work-in-progress area.

Parts of this chapter were published in:

- Ikbal Chammakhi Msadaa, Daniel Câmara, and Fethi Filali. Scheduling and CAC in IEEE 802.16 Fixed BWNs : a Comprehensive Survey and Taxonomy. "IEEE Communications Surveys & Tutorials", 99, Oct. 2010.
- Tijani Chahed, Ikbal Chammakhi Msadaa, Rachid Elazouzi, Fethi Filali, Salah-Eddine Elayoubi, Benoit Fourestié, Thierry Peyre, and Chadi Tarhini. WiMAX Network Capacity and Radio Resource Management. Book chapter in "Radio Resources Management in WiMAX : From theoretical capacity to system simulations", ISBN: 9781848210691, Feb. 2009.

#### **Chapter 5: Adaptive Scheduling with Max-Min Fairness Admission Control**

Despite including the possibility of QoS support, 802.16 MAC protocol does not include a complete solution to offer QoS guarantees for various applications: resource management and scheduling still remain as open issues. In this chapter, we propose a new QoS architecture for PMP 802.16 systems operating in TDD mode over WirelessMAN-OFDM physical layer. It includes a CAC policy and a hierarchical scheduling algorithm. The proposed CAC policy adopts a Min-Max fairness approach making efficient and fair use of the available resources. The proposed scheduling algorithm flexibly adjusts uplink and downlink bandwidth to serve unbalanced traffic. This adaptive per-frame uplink/downlink allocation procedure takes into account the link adaptation capability supported by WiMAX and the data rate constraints of the different types of services. Through simulation, we reveal the efficiency of the proposed CAC scheme and show that our scheduling algorithm can meet the data rate requirements of the scheduling services specified by the IEEE 802.16 Standard.

Parts of this chapter were published in:

- Ikbal Chammakhi Msadaa, Fethi Filali, and Farouk Kamoun. An adaptive QoS Architecture for IEEE 802.16 Wireless Broadband Networks. In MASS 2007, 4th IEEE International Conference on Mobile Ad-hoc and Sensor Systems, Oct. 2007.

#### **Chapter 6: mCoSS: a multi-Constraints Scheduling Strategy for WiMAX Networks**

In this chapter, we propose a multi-Constraints Scheduling Strategy (mCoSS) which maximizes the quality of service (QoS) degree of satisfaction for both real-time and non-real-time traffic in terms of delay and throughput. mCoSS addresses two constraints that were not considered in the QoS solution presented in Chapter 5: latency requirements of real-time connections and support of bursty traffics. In the scheduling strategy presented in this chapter, the access to the network is regulated via a traffic shaper that is inspired from the dual token-buckets shaping mechanism which allows traffic burstiness while protecting contract-conforming connections from misbehaving ones. The modified dual-bucket mechanism is combined with a two-rounds scheduling algorithm reflecting the two levels of service to be expected by each connection. In the first round,

---

the minimum reserved traffic rate and delay constraints are met while in the second round, fairness among flows is ensured over the remaining bandwidth using a weighted fair queuing (WFQ) mechanism. The bandwidth request and grant policy adopted in the proposed strategy takes advantage of the different mechanisms specified by the IEEE 802.16e standard and adapts the choice of the appropriate technique to the service flow QoS constraints and to the current availability of radio resources. Other concerns such as supporting the link adaptation capability and avoiding starvation of best effort traffic are also addressed in the proposed solution.

## **Chapter 7: Mobile WiMAX: a V2I Communications Medium**

Intelligent Transportation Systems (ITS) have been under development since the 80's as part of a global strategy for solving many of our modern life transportation problems. These systems enable people to reach their destinations in a safe, efficient, and comfortable way. In order to reach that goal, several radio access technologies (RAT) such as UMTS, WiFi, WiMAX and 5.9 GHz have been proposed for next generation ITS.

In addition to the 5.9 GHz, which is dedicated to vehicular ad hoc networks networks, mobile WiMAX is expected to play a major role in ITS since it is the only mobile broadband technology currently in use.

In this chapter, we compare mobile WiMAX (based on IEEE 802.16e standard) and 5.9 GHz technology (based on the upcoming IEEE 802.11p standard). We investigate, through simulation, the potential and limitations of both technologies as a communication media for vehicle-to-infrastructure (V2I) communications. The performance of the two systems is evaluated for different vehicle speeds, traffic data rates, and network deployments.

Parts of this chapter were published in:

- Ikbal Chammakhi Msadaa, Pasquale Cataldi, and Fethi Filali. A Comparative Study between 802.11p and Mobile WiMAX-based V2I Communication Networks. In NGMAST 2010, 4th International Conference on Next Generation Mobile Applications, Services and Technologies, July 2010.

## **Appendix A: Topics Related to Mobility Management in WiMAX Networks**

The WiMAX forum estimates that more than 133 million of people will be using the WiMAX technology by the year 2012. From these users, more than 70% are expected to be using the mobile implementation of the technology. From this perspective, mobility management is a key aspect to provide access for these potential 70% of WiMAX users.

This appendix focuses on the latter topic. It describes some concepts and mechanisms introduced by the IEEE 802.16e standard—the amendment of the IEEE 802.16d-2004 standard—which provides enhancements mainly related to mobility management. We also cover the main topics related to WiMAX networks from a mobility perspective and point out the research issues where there is room for contribution.

Parts of this appendix were published in:

- Ikbal Chammakhi Msadaa, Daniel Câmara, and Fethi Filali. Mobility Management in WiMAX Networks. Book chapter in "WiMAX Security and Quality of Service : An End-to-End Perspective". ISBN : 978-0-470-72197-1. Seok-Yee Tang and Peter Muller and Hamid Sharif Ed., July 2010.
-

# Chapter 1

## An overview of WiMAX

In this chapter, we first go through the standardization process of the IEEE 802.16 family of standards. In Section 1.2, we describe the different physical interfaces targeted by the IEEE Std 802.16 as well as the frequency bands for which they have been specified. An overview of the PHY layer is provided in Section 1.3 with a particular insight into the adaptive modulation and coding capability supported by WiMAX. The WiMAX frame formats will be presented in the same section. A brief overview of the MAC layer is given in Section 1.4. This last section only introduces the core functionality of the MAC layer. All the features related to QoS support at MAC level will be further discussed in Chapter 3.

### 1.1 IEEE 802.16/WiMAX standardization

The IEEE 802.16 broadband wireless access (BWA) standard has been developed by the IEEE 802.16 working group (WG) since 1999. The standard was initially designed to support fixed BWA in line-of-sight (LOS) environment in the 10-66 GHz band. It has then been extended to the non-LOS (NLOS) environment in the 2-11 GHz band with the publication of the IEEE 802.16a standard. The IEEE 802.16 task group d (TGd) was later organized to revise and consolidate these standards in a final version, IEEE 802.16-2004 [1], which was approved in 2004.

In December 2005, an amendment of this version: IEEE 802.16e-2005 [2] was published, extending the scope of the standard from fixed to both fixed and mobile environments. This amendment, developed by the IEEE 802.16 TGe, provides enhancements to IEEE Std 802.16-2004 to support subscriber stations moving at vehicular speeds. Both standards have later been rolled up, along with other standards (e.g. 802.16g-2007 related to the Management Plane Procedures and Services), in the IEEE 802.16-2009 standard document [37].

Like for any other technology, the 802.16 standards define a huge set of design alternatives and optional features in order to accommodate the needs of different environments. However, for seek of compatibility between vendor products, only a limited set of mechanisms and certification profiles have been retained by the Worldwide Interoperability for Microwave Access (WiMAX) forum.

Established in 2001, the WiMAX Forum is the entity in charge of promoting and certifying wireless broadband equipments based on the IEEE 802.16 and the European telecommunications standards institute (ETSI) HiperMAN standards. Moreover, it was the WiMAX forum that commercialized the 802.16 family of standards, officially called WirelessMAN in IEEE, under the name "WiMAX". For the rest of this thesis, the terms WiMAX and IEEE 802.16 will be used inter-

---

changeably.

## 1.2 IEEE 802.16 frequency spectrum and PHY interfaces

As mentioned before, the IEEE 802.16 standard specifies the air interface for BWA systems in two different bands: 10-66 GHz and sub 11 GHz. Due to short wavelength, the 10-66 GHz band provides a physical environment where LOS transmission is required and where multipath effect is negligible [38]. The channels used in this band are large: typically 25 or 28 MHz. The PHY interface dedicated to this band is WirelessMAN-SC based on single-carrier (SC) modulation. For the band "below 11 GHz" and more specifically from 2 to 11 GHz, the IEEE 802.16 standard defines two air interfaces: WirelessMAN-OFDM and WirelessMAN-OFDMA based on the orthogonal frequency division multiplex (OFDM) and orthogonal frequency division multiple access (OFDMA) modulations, respectively. These two interfaces operate in 2-11 GHz licensed bands. A third PHY interface WirelessHUMAN (High-speed Unlicensed Metropolitan Area Network) is proposed for 2-11 GHz license-exempt bands. The standard does not specify the modulation technique used in this interface, nevertheless, the unlicensed frequency is included in fixed WiMAX certification. Note that for the 2-11 GHz band, due to longer wavelength, LOS is not required and multipath effect maybe significant. Among the four air interfaces presented in this section, WiMAX only considers WirelessMAN-OFDM and WirelessMAN-OFDMA PHY layers. The two air interfaces are typically dedicated to fixed and mobile systems, respectively.

The common key point between these two PHYs is the use of OFDM which is able to cope with severe channel conditions. Indeed, compared to single-carrier modulation, OFDM offers a higher bandwidth efficiency using a digital multi-carrier modulation. The technique consists in dividing a high data rate stream into several parallel data streams and modulating each of them on a separate subcarrier. These subcarriers are closely-spaced, yet not interfering since they are orthogonal to each other.

## 1.3 An overview of WiMAX PHY

Nodes belonging to the same WiMAX network, share the same wireless medium using one of the two modes specified in the IEEE 802.16 standard: the two-way PMP mode (mandatory) and the mesh mode (optional). The main difference between the two modes is that in mesh mode, subscriber stations (SSs) have the possibility to communicate with each other directly or through the base station (BS), depending on the transmission algorithm in use: distributed, centralized, or a combination of both. In PMP mode however, which is the only mode for sharing media considered in this thesis, a central BS receives and coordinates all the transmissions occurring between SSs. The SSs within a given antenna sector receive the same transmission broadcast by the BS on the downlink channel (DL). Each SS is required to capture and process only the traffic addressed to itself (or to a broadcast or multicast group it is a member of). On the uplink channel (UL) however, the multiple user access is possible either through time division multiple access (TDMA) in WirelessMAN-OFDM or using frequency division multiple access (FDMA) in WirelessMAN-OFDMA, both associated with OFDM modulation technique. Downlink and uplink channels are duplexed using one of the two following techniques: Frequency Division Duplexing (FDD) and Time Division Duplexing (TDD). FDD typically divides the frequency band into two bands: one for the downlink transmission and another one for uplink transmission. In contrast to FDD, TDD systems use the same band for both downlink and uplink and divide the frame, in the time domain, into a DL subframe and an UL subframe. This mode offers channel reciprocity and the possibility

---

Modulation	Coding rate	Receiver SNR (dB)
BPSK	1/2	3.0
QPSK	1/2	6.0
	3/4	8.5
16-QAM	1/2	11.5
	3/4	15.0
64-QAM	2/3	19.0
	3/4	21.0

Table 1.1: Receiver SNR assumptions (WirelessMAN-OFDM) - Table 312 - [37]

of adapting the DL/UL ratio in a dynamic and asymmetric way. It is worth mentioning that only the TDD mode of WiMAX has been accepted as an IMT2000 technology in 2007. Nevertheless, in order to offer more flexibility in the channel bandwidth options, the WiMAX Forum has put a special effort in adding FDD-specific part [39] to Release 1.5 which was approved in August 2009.

In addition to TDD and FDD duplex, OFDM modulation, and support of time and frequency multiple access techniques, WiMAX systems support the use of adaptive modulation and coding (AMC) in order to combat and even take advantage of the channel state fluctuations encountered in wireless propagation environments. We dedicate the next section to describing the latter technique.

### 1.3.1 Link adaptation, modulation, and coding

The adaptive modulation and coding (AMC), also referred to as link adaptation capability, is a powerful technique used by WiMAX technology to strengthen the robustness of the communication to the highly varying channel conditions. This is achieved by employing a robust modulation and coding scheme (MCS) i.e. transmitting at low data rates when the channel is poor and increasing the data rate i.e. using a more efficient MCS when the channel is good. The modulation techniques supported by WiMAX technology are: BPSK, QPSK, 16-QAM, and 64-QAM. For channel coding, three different forward error correction (FEC) types are supported by WiMAX: convolutional codes, turbo codes, and block codes. These channel coding techniques are used to add to the information bits redundant bits which are intended to increase the coding gain and correct the bit errors occurred during transmission. Combined with the different modulation and coding rates proposed by the IEEE 802.16 standard, these FEC types lead to 52 possible configurations called "burst profiles".

The mechanism used to choose the most appropriate per-frame and per-user MCS, and manage the DL and UL burst profiles of each SS, is not fully specified by the standard. Moreover the guidelines and recommended policies depend on the PHY layer in use (OFDM or OFDMA). Nevertheless, the basic idea consists in adapting the choice of the most appropriate burst profile, identified by a DIUC/UIUC (DL/UL interval usage channel), to the channel SNR (signal-to-noise ratio) measured at the receiver. Tables 1.1 and 1.2 report the MCSs recommended by the IEEE 802.16 standard [37] for given values of the receiver SNR in WirelessMAN-OFDM and WirelessMAN-OFDMA, respectively. Note that these are only order of magnitudes of SNRs obtained for specific requirements (a BER of  $10^{-6}$  measured after FEC) and channel conditions<sup>1</sup>.

<sup>1</sup>The reported values are derived in an AWGN environment. Table 1.1 SNRs assume the use of Reed-Solomon convolutional coding (RS-CC) while Table 1.2 SNRs are obtained assuming the use of a tail-biting convolutional code.

Modulation	Coding rate	Receiver SNR (dB)
QPSK	1/2	5.0
	3/4	8.0
16-QAM	1/2	10.5
	3/4	14.0
64-QAM	1/2	16.0
	2/3	18.0
	3/4	20.0

Table 1.2: Receiver SNR assumptions (WirelessMAN-OFDMA) - Table 545 - [37]

### 1.3.2 WirelessMAN-OFDM

In the IEEE 802.16, the channel consists of fixed-length frames, as shown in Figure 1.1. Each frame is divided into DL and UL subframes. [1] specifies that, when using TDD, the UL subframe and DL subframe durations shall vary within the same shared frame. The downlink subframe consists of one single PHY PDU while the uplink subframe consists of two contention intervals followed by multiple PHY PDUs, each transmitted by a different SS. The first contention interval is used for ranging which is the process of adjusting the radio frequency (RF). The second interval may be used by the SSs to request bandwidth since bandwidth is granted to SSs on demand. Two gaps separate the downlink and uplink subframes: transmit/receive transition gap (TTG) and receive/transmit transition gap (RTG). These gaps allow the BS to switch from the transmit to receive mode and vice versa.

The downlink PHY PDU consists of one or more bursts, each transmitted with a specific burst profile. A burst profile is a set of parameters describing the transmission properties (modulation type, forward error correction (FEC) type, etc.) corresponding to an interval usage code (IUC). Each SS is required to adapt the IUC in use (a DIUC for the downlink and an UIUC for the uplink) based on measurements on the physical layer. The length of each burst is set by the BS. Indeed, at the beginning of each frame, the BS schedules the uplink and downlink grants (by mechanisms that are outside the scope of the standard [1, 2]) and then broadcasts the downlink frame prefix (DLFP), the DL-MAP and the UL-MAP informing the SSs of its scheduling decisions. The DLFP describes the location and profile of the first downlink bursts (at most four). SSs using the same DIUC are advertised as a single burst. The DL-MAP, when sent, describes the location and profile of the other downlink bursts—if they exist. However, the IEEE 802.16 standard specifies that, at least one full DL-MAP must be broadcast within the Lost DL-MAP Interval even if there are less than five bursts. The UL-MAP should be transmitted in each frame. It contains information elements (IE) that indicate the types and the boundaries of the uplink allocations directed to the SSs.

The profile of each downlink and uplink burst are specified in the downlink channel descriptor (DCD) and uplink channel descriptor (UCD), respectively. The BS broadcasts the DCD and the UCD messages periodically—every DCD/UCD Interval—in order to define the characteristics of the downlink and uplink physical channels. Referring to Figure 1.1, we note that each burst consists of one or more MAC PDUs. Each MAC PDU begins with a fixed-length MAC header followed by a payload and a cyclic redundancy check (CRC) field. The burst may also contain padding bytes since each burst must consist of an integer number of OFDM symbols. UL bursts begin with a preamble used for PHY synchronization.

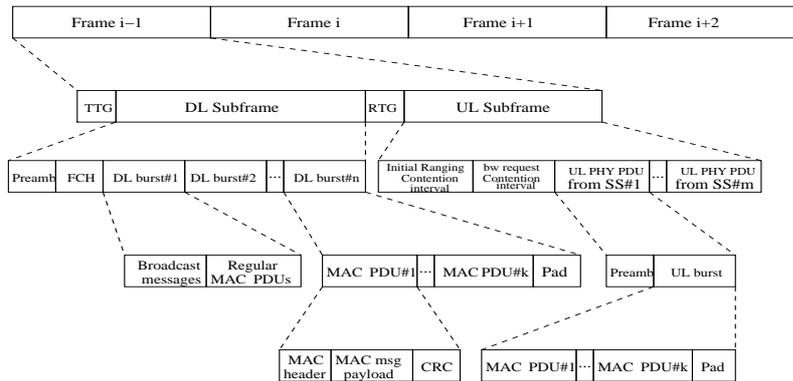


Figure 1.1: OFDM Frame Structure with TDD

### 1.3.3 WirelessMAN-OFDMA

#### 1.3.3.1 Subchannelization

WirelessMAN-OFDMA supports channel bandwidths of 3.5 to 20 MHz with 128, 512, 1024, or 2048 subcarriers. The available subcarriers are grouped into groups of subcarriers called subchannels. This subchannels might be formed using either distributed or adjacent subcarriers which correspond to the two modes of subcarriers permutation supported by the IEEE 802.16 standard. In distributed permutation, also called diversity permutation, subcarriers forming the same subchannel are pseudorandomly taken over a the frequency spectrum, thus achieving frequency diversity gain. In adjacent permutation however, the subcarriers belonging to a same subchannel are physically adjacent. This permutation is called band adaptive modulation and coding (band AMC). This technique enhances the spectrum efficiency by selecting the user that have a strong channel and by choosing the MCS that maximizes the band efficiency, thus achieving multi-user diversity through frequency-selective resource allocation.

#### 1.3.3.2 Band AMC Permutation mode

The basic allocation unit in an OFDMA system depends on the considered subchannelization mode and could vary from DL to UL. In this section, we present the specifications related to band AMC subchannelization. In the latter mode, the basic allocation unit in DL and UL is called a bin. One bin, as shown in Figure 1.3, corresponds to nine contiguous subcarriers, one pilot and eight data subcarriers, within an OFDMA symbol. A group of four adjacent bins in the frequency domain is called a physical band, and a grouping of physical bands forms a logical band. The smallest time-frequency resource that can be allocated is called a "slot" and corresponds to six contiguous bins within the same logical band. Thus, as illustrated in Figure 1.3, a slot consists of one bin over six symbols, two bins over three symbols, or three bins by two symbols. It could also correspond to a default type formed by six contiguous bins enumerated as follows: starting from the lowest bin in the first symbol to the last bin in the next symbol. This enumeration is also applied for the other three types, as shown in Figure 1.3. In all cases, each slot consists of 48 data subcarriers within a subchannel.

#### 1.3.3.3 OFDMA frame structuring

In contrast to OFDM, a second dimension corresponding to subcarriers is introduced in OFDMA. Nevertheless, the frame structures in OFDM and OFDMA are quite similar. Figure 1.2 shows

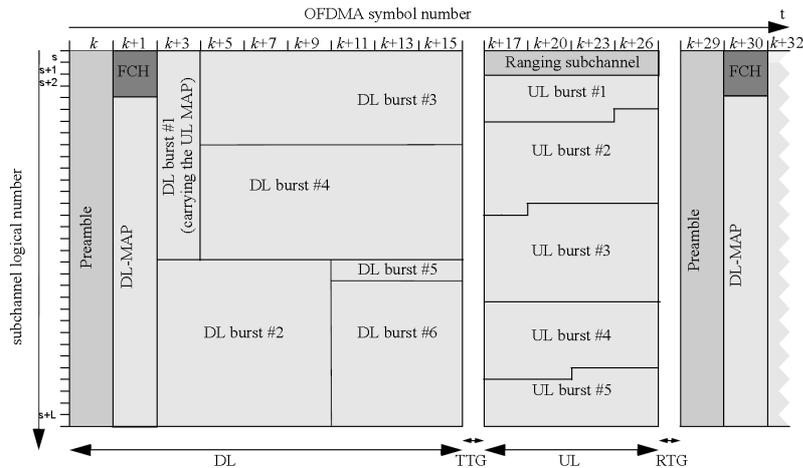


Figure 1.2: Example of OFDMA frame in TDD

an example of an OFDMA frame structure in TDD mode. The frame is divided into a DL subframe and an UL subframe. The DL subframe starts with a preamble followed by a 24-bit FCH containing the DLFP which specifies the length and the repetition coding used for the DL-MAP message. Then the DL-MAP and UL-MAP messages are broadcast. They include the duration and the profiles of the downlink and uplink bursts, respectively. In addition to the data regions allocated for the users, the UL subframe consists of portions reserved for contention-based access. These portions are mainly used for bandwidth requests, initial, periodic and handover ranging or to give the MSs the opportunity to acknowledge (ACK/NACK) DL transmissions. Part of these contention-access slots might be used by the BS to allocate a channel quality information channel (CQICH) for the MSs to transmit periodic CINR reports. Like in OFDM, a TTG and RTG gaps are inserted between the DL and UL subframes and at the end of the frame allowing the BS to switch from transmitting to receiving mode and vice-versa.

## 1.4 An overview of WiMAX MAC

### 1.4.1 A connection-oriented MAC

In addition to the physical layer, the standard defines a connection-oriented MAC layer where all the data transmissions occur within the context of a unidirectional transport connection. Each connection, identified by a unique 16-bit connection ID (CID), is associated to a service flow (SF) whose characteristics provide the QoS requirements to apply for the protocol data units (PDUs) exchanged on that connection. In addition to transport connections dedicated to data transmissions, each SS is assigned at the initialization two pairs of management connections, basic connections (DL and UL) and primary management connection (DL and UL). An optional third pair of secondary management connections might also be established. The use of these three pairs of connections reflects the three different levels of QoS associated to the different management messages exchanged between the BS and the SS:

- Basic connection: used to transfer short and time-critical MAC management messages such as the messages reporting the channel measurements (REP-REQ and REP-RSP).
- Primary management connection: used to transfer longer and more delay-tolerant messages

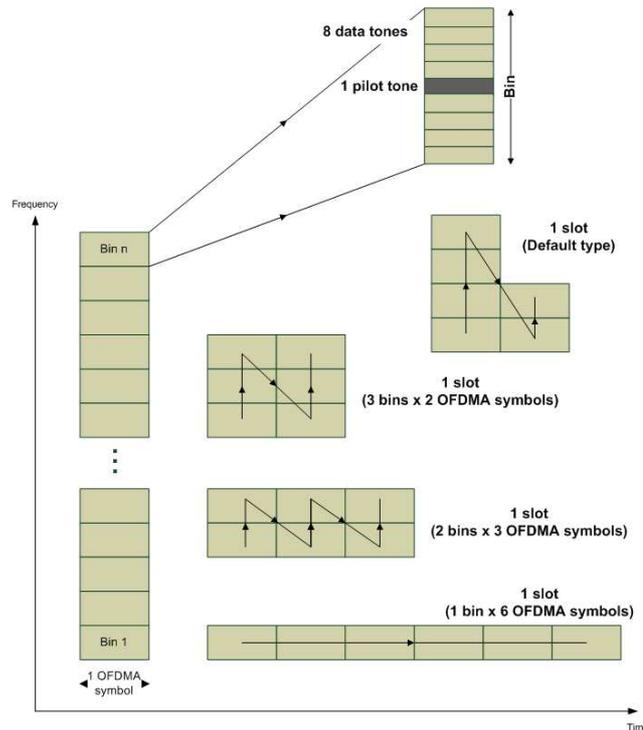


Figure 1.3: Slot structure in Band AMC Permutation

like those used to create a new service flow e.g. dynamic service addition request (DSA-REQ).

- Secondary management connection: this third pair is required only for managed SSs and is used to transfer delay-tolerant, standard-based messages e.g. dynamic host configuration protocol (DHCP), simple network management protocol (SNMP) messages.

## 1.4.2 IEEE 802.16 protocol stack

Figure 1.4 illustrates the IEEE 802.16 standard reference model for the data plane. As shown in this figure, the MAC layer specified by the standard consists of three sublayers:

- a service-specific convergence sublayer (CS): the receiving CS accepts the MAC SDUs from the peer MAC SAP (service access point) and delivers them to the upper layers. At the transmitting entity, the CS is responsible for delivering, to the MAC SAP, the upper layer PDUs received through the CS SAP. A classification process is performed at this level, based on a set of protocol-specific matching criteria, in order to associate the network SDUs to the proper MAC service flow identifier (SFID). This process facilitates the MAC SDUs delivery with the appropriate QoS constraints. The CS sublayer may also include additional functions such as payload header suppression (PHS).
- a MAC common part sublayer (CPS) providing the core functionality of the MAC layer: network-entry process, connection establishment and maintenance, control and signalling, bandwidth allocation, and QoS support.

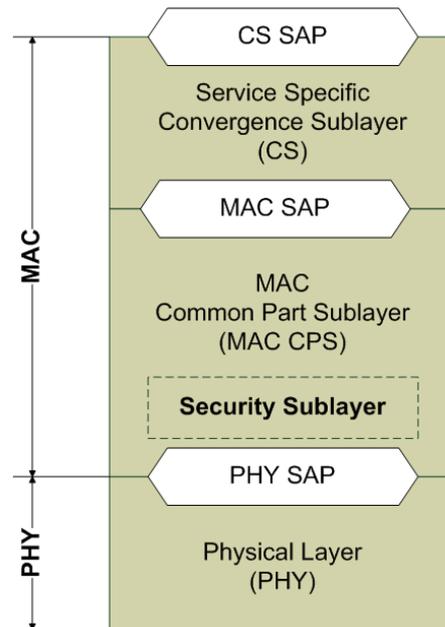


Figure 1.4: IEEE Std 802.16 Data Plane Protocol Reference Model

- a separate security sublayer providing authentication and privacy to MSs and protecting the operators from theft of service. This is performed through the use of (i) an encapsulation protocol for securing packet data across the BWA network and (ii) a key management protocol (PKM) securing the securing data from the BS to the MS.

### 1.4.3 MAC PDUs formats, construction and transmission

As shown in Figure 1.1, the basic structure of a MAC PDU consists of a 48-bit generic MAC header, optionally followed by a payload and a CRC field. The use of CRC is mandatory for both OFDM and OFDMA PHYs. For DL traffic, only the generic MAC header is defined whereas for UL, the standard defines two types of headers. The first one is the generic header whose format is shown in Figure 1.5(a). It begins each MAC PDU containing either a MAC management message or CS data. In this header format, the header type (HT) field is set to 0. The second header format, where HT is set to 1, is a signalling header that is not followed by any MAC PDU payload or CRC. This second type is dedicated to the transfer of short signalling information such as bandwidth request or feedback information (e.g. the UL Tx Power report header and the CINR report header) which does not require the overhead associated to a payload [38]. Figure 1.5 (b) shows an example of this signalling header used for bandwidth request. The BR field in this header refers to the number of UL bytes requested by the SS and should be independent of the MCS in use. More details about the use of the bandwidth request header will be given in Section 3.1.3 when we introduce the bandwidth request mechanisms proposed by the IEEE Std 802.16.

In the MAC PDU, with the generic MAC header, there exist different per-PDU subheaders like the fragmentation, or the grant management (GM) subheader and one per-SDU subheader which is the packing subheader. Packing and fragmentation, along with the concatenation techniques are used by the MAC protocol to enhance the efficiency of the air interface. Indeed, a long MAC SDU (or long MAC management message) may be divided into multiple MAC PDUs and multiple MAC SDUs (or MAC management messages) may be combined into a single PDU if their lengths

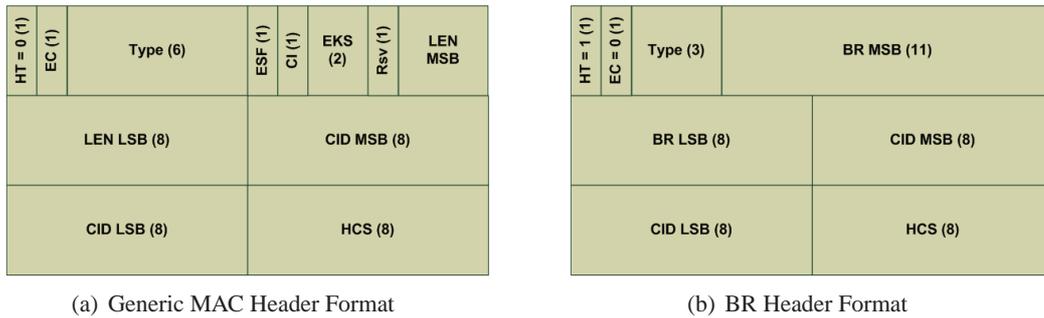


Figure 1.5: IEEE 802.16 MAC Headers Formats

are short. The former process is called "fragmentation" and the latter "packing". The process of combining multiple PDUs into a single burst like we have seen in Section 1.3.2 when describing the OFDM TDD frame format is called "concatenation".

## 1.5 Conclusion

This chapter is aimed at providing an overview of the main features supported by PHY and MAC layers specified by the IEEE 802.16 standard. Therefore, we have first gone through the standardization process of the IEEE 802.16 family of standards. Then, we have described both WirelessMAN-OFDM and WirelessMAN-OFDMA air interfaces targeted by the IEEE Std 802.16, the frequency bands for which they have been specified as well as their corresponding frame formats. An overview of the PHY layer has been provided in this chapter with a particular insight into the adaptive modulation and coding capability supported by WiMAX. As for MAC layer, only the core functionality, necessary to the understanding of the performance study carried out in Chapter 2, was described in this chapter. All the features related to QoS support at MAC level will be further discussed in Chapter 3. Indeed, because there are so many concepts to be introduced in this context, we have preferred to dedicate a whole chapter to this purpose.



---

## Chapter 2

# Performance Analysis of OFDM-based WiMAX Networks

In this chapter, we evaluate the performance bounds of WiMAX systems under different physical and MAC parameters settings. The saturation throughput that can be reached in 802.16 networks is investigated through several scenarios in which we vary for instance the frame duration, the channel bandwidth, and the modulation and coding scheme (MCS) in use. An analytical framework was developed based on technical properties and system profiles specified by the IEEE 802.16 standard for systems using the WirelessMAN-OFDM air interface. The obtained results outline the importance of considering the MAC and physical overhead when evaluating the performance of 802.16 networks. They also highlight the impact of packing and fragmentation techniques, proposed by IEEE 802.16 standard, on the MAC performance and show the trade-off to make between decreasing the channel bandwidth and increasing the resulting saturation throughput. The remainder of this chapter is structured as follows. An analytical framework considering technical properties of WirelessMAN-OFDM PHY variant is developed in Section 2.1. The performance evaluation study is detailed in Section 2.2. Section 2.3 concludes the chapter by outlining the main obtained results.

### 2.1 Analytical framework

In this section, we first need to detail some technical features related to WirelessMAN-OFDM PHY and that were not mentioned in Chapter 1. Secondly we carry out an analytical study of the OFDM PHY frame structure described in Section 1.2. This study is aimed at giving analytical expressions of the saturation throughput that may be reached in 802.16 networks while taking into account the MAC and PHY overhead. As mentioned in Section 1.2, WirelessMAN-OFDM PHY is designed for frequencies below 11 GHz where LOS is not necessary and where multipath may be significant. To collect multipath, a cyclic prefix (CP) is used. As depicted in Figure 2.1(a), this prefix corresponds to a copy of the last  $T_g$  of the useful symbol time  $T_b$  of an OFDM symbol  $T_{sym}$ . The OFDM symbol transmission time is then expressed as follows:  $T_{sym} = T_g + T_b$ ; where the guard time  $T_g$  is given by:  $T_g = g * T_b$ .  $g$  corresponds to the ratio of CP time to useful time. The possible values of  $g$  are: 1/4, 1/8, 1/16, and 1/32 [2].

As for the frequency domain structure, an OFDM symbol, described by Figure 2.1.b, is composed of data subcarriers (for data transmission), pilot subcarriers (for estimation purposes) and null subcarriers such as guard subcarriers. The total number of subcarriers corresponds to the fast Fourier transform (FFT) size  $N_{fft}$ . According to [2],  $N_{fft} = 256$ . Let  $BW$ ,  $n$  and  $F_s$  denote the

---

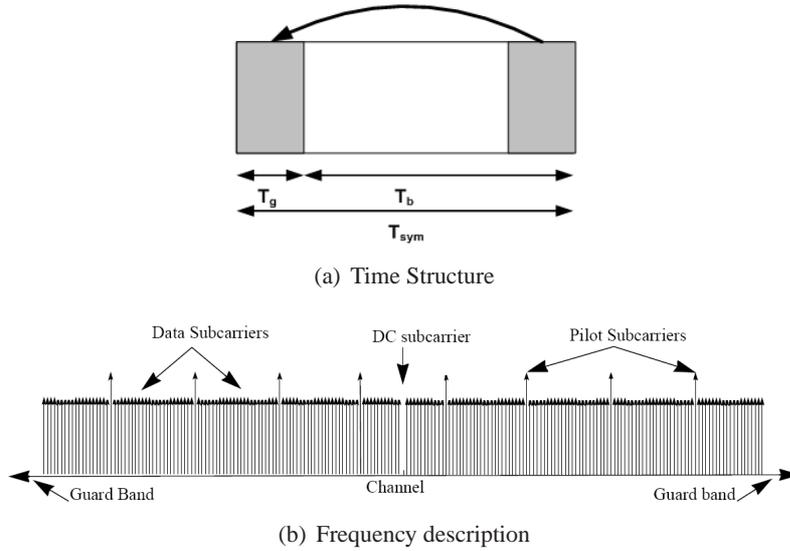


Figure 2.1: OFDM Symbol structure

nominal channel bandwidth, the sampling factor and the sampling frequency, respectively. The sampling frequency corresponds to:  $F_s = n * BW$ . The value of the sampling factor  $n$  depends on the channel bandwidth  $BW$  as it is illustrated by Table 2.1. The possible values of  $BW$  correspond to those specified in the system profiles proposed by the IEEE 802.16 standard [2] for systems operating with the WirelessMAN-OFDM air interface. As shown in Table 2.1, five PHY profiles are specified for these systems, each corresponding to a channel bandwidth. Suppose that  $\Delta f$  stands for the subcarrier spacing, then:  $\Delta f = F_s/N_{fft}$  and the useful time is given by:  $T_b = 1/\Delta f$ .

For a given system configuration ( $BW$  and  $g$  fixed), the duration of an OFDM symbol is fixed. However, in terms of data, the number of information bits per OFDM symbol varies depending on the modulation and coding scheme (MCS) in use. Indeed, if the selected MCS has a constellation size  $M$  with efficiency  $k_{MCS} = \text{Log}_2(M)$  and a coding rate  $CR_{MCS}$  the number of information bits per symbol is computed as follows.

$$N_{MCS}^{bpsym} = N_{data-sub} \times k_{MCS} \times CR_{MCS} - 8 \quad (2.1)$$

where:

- $N_{data-sub}$  stands for the number of data subcarriers ( $N_{data-sub} = 192$ ).
- The "-8" refers to the 0x00 tail byte at the end of each OFDM symbol.

For 16-QAM 3/4, for instance,  $N_{16QAM-3/4}^{bpsym} = 192 * 4 * 3/4 - 8 = 568$ .

Let us consider the OFDM PHY structure illustrated in Figure 2.2. First we focus on fixed-size fields/intervals. Therefore let us consider the following parameters:

- $T_{frame}$ : duration of a time frame (in seconds).
- $T_{av}$ : time duration (in seconds), still available in the frame. Initially, we have:  $T_{av} = T_{frame}$ .

System Profile Identifier	Channel Bandwidth $BW$ (MHz)	Sampling factor $n$
profP3_1.75	1.75	8/7
profP3_3	3	86/75
profP3_3.5	3.5	8/7
profP3_5.5	5.5	316/275
profP3_7	7	8/7

Table 2.1: WirelessMAN-OFDM System Profiles

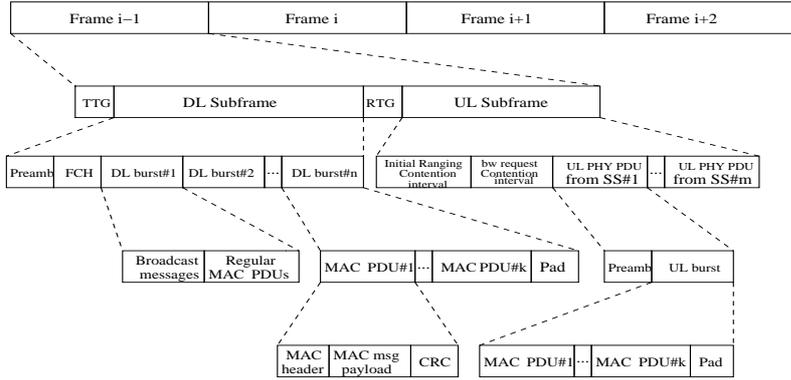


Figure 2.2: OFDM Frame Structure with TDD

- $T_{sym}$ : duration of an OFDM symbol (in seconds).
- $T_{pream}^{(short)}$ : duration of a short preamble (in seconds),  $T_{pream}^{(short)} = T_{sym}$  [1].
- $T_{pream}^{(long)}$ : duration of a long preamble,  $T_{pream}^{(long)} = 2 * T_{sym}$  [1].
- $T_{ttg}$ : duration of a transmit/receive transmission gap (in seconds).
- $T_{rtg}$ : duration of a receive/transmit transmission gap (in seconds).
- $T_{dlfp}$ : duration of the DLFP (in seconds). [1] specifies that  $T_{dlfp} = T_{sym}$ .
- $S_{opp}^{(rng)}$  and  $S_{opp}^{(bw)}$ : size (in OFDM symbols) of a ranging and bandwidth request opportunity, respectively.
- $N_{opp}^{(rng)}$  and  $N_{opp}^{(bw)}$ : number of ranging and bandwidth request opportunities during the contention interval, respectively.
- $T_{opp}^{(rng)}$  and  $T_{opp}^{(bw)}$ : duration (in seconds) of a contention ranging and bandwidth request interval, respectively, expressed as follows:

$$T_{opp}^{(rng)} = N_{opp}^{(rng)} * S_{opp}^{(rng)} * T_{sym} \quad (2.2)$$

$$T_{opp}^{(bw)} = N_{opp}^{(bw)} * S_{opp}^{(bw)} * T_{sym} \quad (2.3)$$

Note that all the above-cited parameters are multiples of  $T_{sym}$ , so we can deduct their respective durations from  $T_{av}$  since  $T_{av}$  should always be kept as an integer number of OFDM symbol duration.

$$T_{av} = T_{frame} - \left( T_{pream}^{(long)} + T_{dlfp} + T_{ttg} + T_{opp}^{(rng)} + T_{opp}^{(bw)} + T_{rtg} \right) \quad (2.4)$$

Recall that the first DL burst contains the broadcast MAC control messages: DCD, UCD, DL-MAP, and UL-MAP. The sizes of these messages depend on the number of DL/UL burst profiles described in the DCD/UCD messages, and on the number of DL/UL IEs specified in the DL-MAP/UL-MAP messages, respectively. However since we are interested in the performance bounds of 802.16 systems, we will consider only one SS and one BS<sup>1</sup>. We also assume that the SS sends continuously to the BS and does not receive any data from it. It is important to mention that a descriptor should be included into DCD message for each DIUC used in the DL-MAP except those associated with Gap, End of Map and Extended IEs. Thus since we assume that no data is transmitted on the downlink, only one DL burst profile is needed to describe the transmission properties of the first DL burst carrying MAC management messages. As for the UL, a burst descriptor shall be included into the UCD message for each UIUC that is to be used in the UL-MAP. Yet, in addition to the end of map IE and to the data grant IE that will specify the amount of bandwidth granted to the SS, an initial ranging IE, and a request IE should be specified in the UL-MAP message to draw the limits of the initial ranging and bandwidth request contention intervals. In our case, these two intervals will be reserved to the single SS belonging to the network. Each of these four IEs will be associated to an UIUC.

Based on the above considerations, let us define the following parameters:

- $S_{dcd}$ : size (in bytes) of a DCD message specifying one DL burst profile.
- $S_{ucd}$ : size (in bytes) of a UCD message specifying four UL burst profiles.
- $S_{dlmap}$ : size (in bytes) of a DL-MAP message that does not specify any burst: it corresponds to the minimum size of a DL-MAP—containing only an end of map IE. Since we have only one DL burst, its limits are specified in the DLFP.
- $S_{ulmap}$ : size (in bytes) of an UL-MAP message containing four IEs—data grant, initial ranging, request, and end of map IE.

These sizes are computed with respect to the type-length-value (TLV) encoding form specified by the standard [1]. They include the MAC overhead (generic header and CRC field). Since DCD, UCD and DL-MAP messages are sent periodically, let  $Send_{dcd}$ ,  $Send_{ucd}$ , and  $Send_{dlmap}$  denote three boolean variables indicating whether a DCD, an UCD or a DL-MAP message will be sent in the current  $i^{th}$  frame, respectively. These parameters are set to 1 each time the timers associated to the following intervals expire: DCD Interval, UCD Interval, and Lost DL-MAP Interval, respectively. As for UL-MAP message, it must necessarily exist in each frame.

To compute the length of the first DL burst, we should take into account the possibility of padding since every burst should consist of an integer number of OFDM symbols. This rule is to be respected each time a burst size is updated. Given a burst  $k$  and its modulation and coding scheme, the number of padding bits is computed such that:

$$\frac{L_{bst}[k] + L_{pad}[k]}{L_{sym}[k]} = n; \quad n \in \mathbb{N} \quad \text{and} \quad L_{pad}[k] < L_{sym}[k] \quad (2.5)$$

where:

---

<sup>1</sup>Just a few modifications are needed to adapt the analytical study to a more general case involving many SSs with different DIUC/UIUC.

- $L_{bst}[k]$  is the number of bits transmitted in burst  $k$  (payload, MAC, and Physical overhead) except the padding bits.
- $L_{pad}[k]$  is the number of padding bits sent in burst  $k$ .
- $L_{sym}[k]$  is the number of bits per OFDM symbol for the burst  $k$ .

Applying (2.5) to the first burst characterized by  $L_{bst}[1]$ ,  $L_{pad}[1]$ , and  $L_{sym}[1]$ , we obtain:

$$L_{bst}[1] = \left( Send_{dcd} * S_{dcd} + Send_{ucd} * S_{ucd} + Send_{dlmap} * S_{dlmap} + S_{ulmap} \right) * 8 \quad (2.6)$$

and  $L_{pad}[1] = compute\_pad(L_{sym}[1], L_{bst}[1])$ ; where  $compute\_pad()$  is a function that returns the number of padding bits necessary for a burst  $k$  given its length and its number of bits per OFDM symbol:

$$compute\_pad(L_{sym}[k], L_{bst}[k]) = L_{sym}[k] - (L_{bst}[k] \% L_{sym}[k]) \quad (2.7)$$

Once  $L_{bst}[1]$  and  $L_{pad}[1]$  are computed, the available time is updated as follows:

$$T_{av} = T_{av} - \frac{L_{bst}[1] + L_{pad}[1]}{L_{sym}[1]} * T_{sym} \quad (2.8)$$

Referring to Figure 2.2, note that all the durations corresponding to MAC management messages, contention intervals, gaps and preambles were considered in the above study. A short preamble duration  $T_{pream}^{(short)}$ —necessary for SS PHY synchronization—should nevertheless be subtracted from the remaining frame duration to get the whole duration available for data transmission:  $T_{av} = T_{av} - T_{pream}^{(short)}$ .

Recall that our main objective is to determine the performance bounds of IEEE 802.16 systems. Therefore it is interesting to compute the maximum number of PDUs  $N_{pdu}^{max}$  that may be transmitted by the SS during the available time. This parameter depends on the considered size of the MAC SDU ( $S_{pkt}$ ), on the modulation and coding scheme used for the UIUC in addition to other PHY parameters like the channel bandwidth  $BW$  and the frame duration  $T_{frame}$ .

$$N_{pdu}^{max}(S_{pkt}) = \frac{(T_{av}/T_{sym}) * L_{sym}[k]}{(S_{gmh} + S_{pkt} + S_{crc}) * 8} \quad (2.9)$$

As can be seen in (2.9), the MAC overhead corresponding to the CRC field and to the MAC generic header and resulting from the transmission of each MAC PDU, are taken into account. Based on (2.9), the maximum MAC goodput—corresponding to the maximum IP throughput (in bps)—that can be reached in such a configuration of 802.16 networks, can be derived as follows:

$$Thput^{max}(S_{pkt}) = \frac{N_{pdu}^{max}(S_{pkt}) * S_{pkt} * 8}{T_{frame}} \quad (2.10)$$

parameter's name	parameter's value
$g$	1/4
$T_{ttg}$	$2 * T_{sym}$
$T_{rtg}$	$2 * T_{sym}$
$N_{opp}^{(rng)}$	1
$S_{opp}^{(rng)}$	3 OFDM symbols
$N_{opp}^{(bw)}$	1
$S_{opp}^{(bw)}$	1 OFDM symbol
$N_{dl\_bp}$	1
$N_{ul\_bp}$	4
$N_{dlmap\_ie}$	1
$N_{ulmap\_ie}$	4

Table 2.2: Physical and MAC parameters

## 2.2 Performance evaluation

As mentioned before, the main parameter investigated in our study is the saturation throughput. The saturation throughput is defined as the highest data rate that could be achieved in the medium. This metric is very important in wireless networks and provides an absolute limit of the amount of data packets that could be successfully sent in the channel. The value of the saturation throughput depends on the overhead induced by the medium access control mechanism. The parameters used in this section are given in Table 2.2.  $N_{dl\_bp}$ ,  $N_{ul\_bp}$ ,  $N_{dlmap\_ie}$ , and  $N_{ulmap\_ie}$  stand for the number of DL burst profiles, UL burst profiles, DL-MAP IEs, and UL-MAP IEs, respectively. Other parameters such as the channel bandwidth, the frame duration, and the MCS will be fixed according to the objective of each scenario. The effect of these parameters on MAC efficiency is investigated in several scenarios.

### 2.2.1 Effect of the frame duration and the MCS

To show the impact of the frame duration and the modulation and coding scheme on the MAC goodput—which also corresponds to the IP throughput—we consider two scenarios. In the first one, we set the frame duration to 20 ms and compute the resulting IP throughput for different modulation and coding schemes. In the second one, we fix the modulation and coding scheme to 64-QAM 3/4 and compute the resulting IP throughput for different frame durations. In both scenarios, the channel bandwidth  $BW$  is set to 7 MHz. Figures 2.3(a) and 2.3(b) depict the IP throughput variation, as a function of MAC SDU size, for scenario 1 and scenario 2, respectively.

As expected, the IP throughput increases with the frame duration as shown in Figure 2.3(b) and, as depicted in Figure 2.3(a), the less robust is the burst profile, the higher is the obtained IP throughput. It is interesting to see that for all the modulation and coding schemes considered in the first scenario, the maximum throughput is reached for nearly the same packet size (more than 100 bytes) and it remains almost the same. However, as can be seen in Figure 2.3(b), a higher fluctuation on MAC goodput can be observed when the frame duration gets shorter. Indeed for a frame duration of 5 ms, the IP throughput fluctuates from almost 9 Mbps to more than 12 Mbps, depending on the packet size; and the bigger is the MAC SDU size, the higher is the fluctuation. This may be explained by the fact that since the fragmentation capability is disabled, in these first scenarios, the possibility that a big packet cannot be transmitted is more likely to happen when the frame duration is short which increases the resulting throughput. Note that the maximum IP

throughput (19.275 Mbps) obtained for a frame duration of 20 ms and 64-QAM 3/4 as modulation and coding scheme, corresponds to the saturation throughput of the considered systems since it uses the biggest channel bandwidth (7 MHz) specified by the system profiles of the IEEE 802.16 standard, the longest possible value of frame duration (20 ms) and the most efficient modulation and coding scheme (64-QAM 3/4).

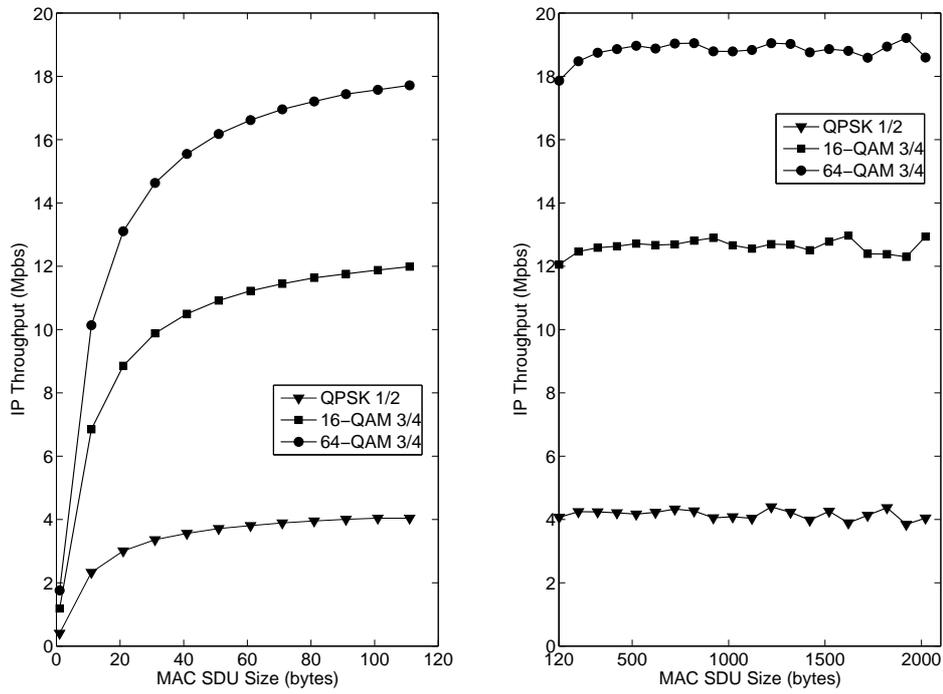
We are still investigating the effect of the frame duration and the MCS on MAC performances. However, in this case, we are more interested on how the whole frame is used: what are the respective proportions—in terms of time—of payload and overhead and what would be the amount of wasted bandwidth in absence of fragmentation. Therefore we introduce two parameters which are the overhead and the wasted time. The overhead (in terms of time) is computed as follows:

$$\begin{aligned}
 Ovhd^{max}(S_{pkt}) = & \quad T_{frame} - T_{av} \\
 + & \quad \left( N_{pdu}^{max}(S_{pkt}) * (S_{gmh} + S_{crc}) * 8 \right. \\
 + & \quad \left. compute\_pad(L_{sym}[k], L_{bst}[k]) \right) \\
 / & \quad L_{sym}[k] * T_{sym}
 \end{aligned} \tag{2.11}$$

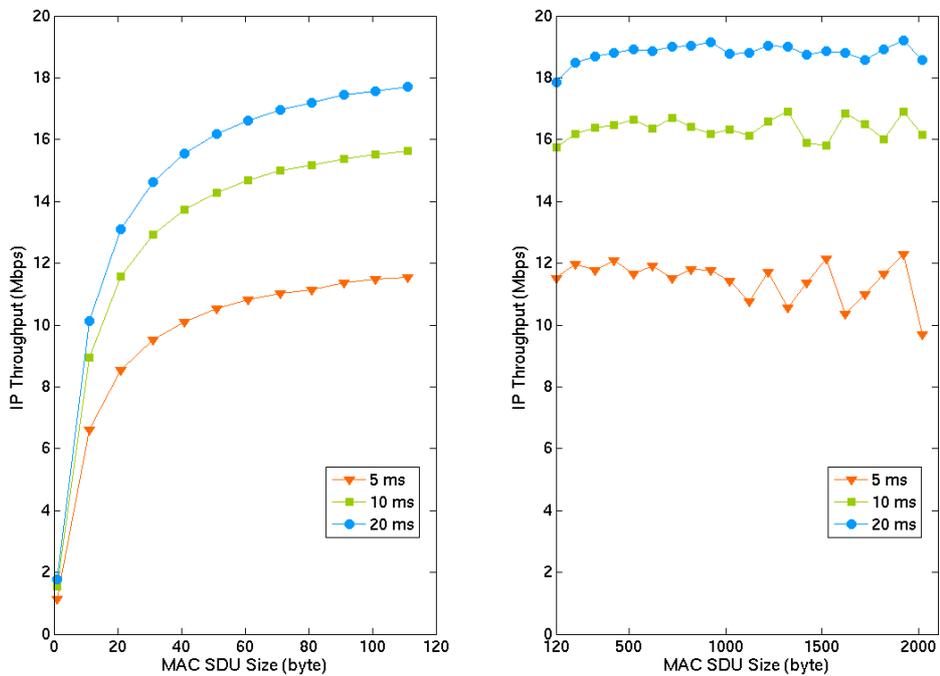
where  $T_{av}$  corresponds to the last value of available time. The overhead corresponds then to the ratio of time—of the frame duration—used for gaps, preambles, contention intervals, and management messages transmission. It also includes the MAC overhead resulting from the transmission of the maximum number of PDUs and the necessary padding. The wasted time corresponds then to the remaining of the frame duration, after omitting the overhead (2.11) and the time needed for the transmission of the maximum number of PDUs (2.9). These three proportions of the time frame are illustrated in Figures 2.4(a) and 2.4(b) for three values of  $T_{frame}$ : 5 ms, 10 ms, and 20 ms. What makes the difference between these two figures is that in Figure 2.4(a), we suppose that the SS uses QPSK 1/2 MCS while in 2.4(b), the use of 64-QAM 3/4 is assumed.

Figure 2.4(a) shows that the longer is the frame, the bigger is the proportion of time reserved for payload transmission and the smaller are the proportions of overhead and wasted time. It is worth mentioning that the overhead may constitutes more than 90% of the frame duration for packets of less than 400 bytes; and this is more likely to happen since almost 75% of the packets of the Internet traffic are smaller than 522 bytes and nearly half of the packets are 40 to 44 bytes in length. In the case of 5 ms frame duration, even for bigger MAC SDUs, the overhead may reach more than 40% of the total frame size.

Now let us compare two frame compositions corresponding to the same frame duration but using two different modulations. If we consider for instance a frame duration of 5 ms in both cases (Figure 2.4(a) and 2.4(b)), we observe that the ratio of overhead increases when using 64-QAM 3/4. This may be explained as follows. Using a less robust modulation (64-QAM 3/4) implies a bigger number of bits per OFDM symbol which offers the possibility of sending more MAC PDUs but also more MAC headers and CRC fields. It also implies the possibility of more padding bits when necessary, in other words more overhead. However having bigger proportion of overhead—in terms of time—does not mean necessarily a decrease of resulting IP throughput since for the same duration, more data can be sent when using 64-QAM 3/4 than when using QPSK 1/2, as we have seen in Figure 2.3(b). Also when comparing Figure 2.4(a) and 2.4(b), we notice that the ratio of wasted time decreases in the case of 64-QAM 3/4, which decreases the effect of absence of fragmentation. Indeed having the possibility of sending more data within the same duration increases the chance of sending even big MAC PDUs and then saving bandwidth.

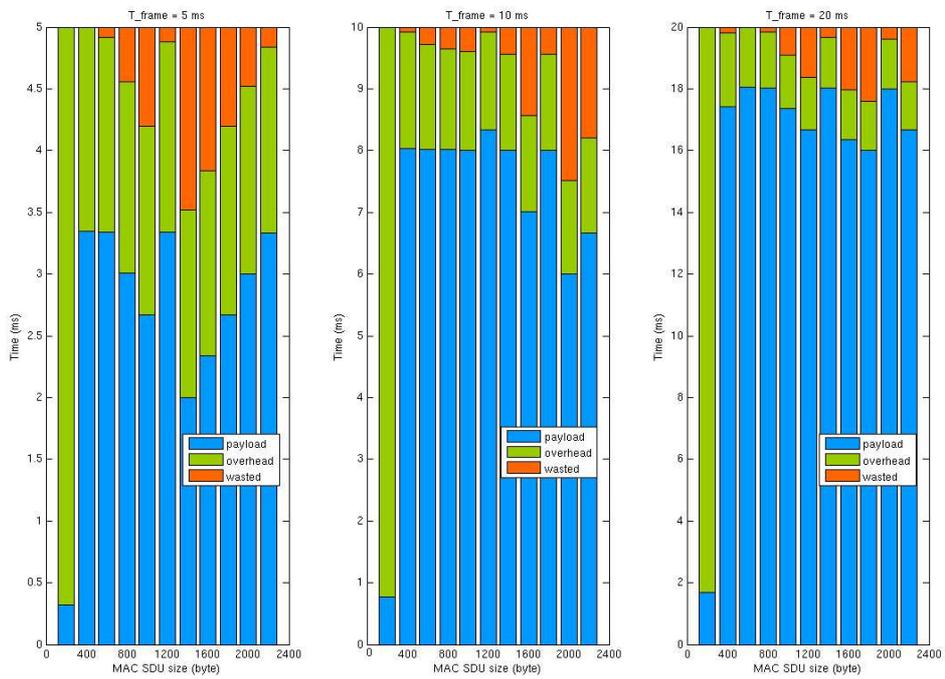


(a) 20 ms

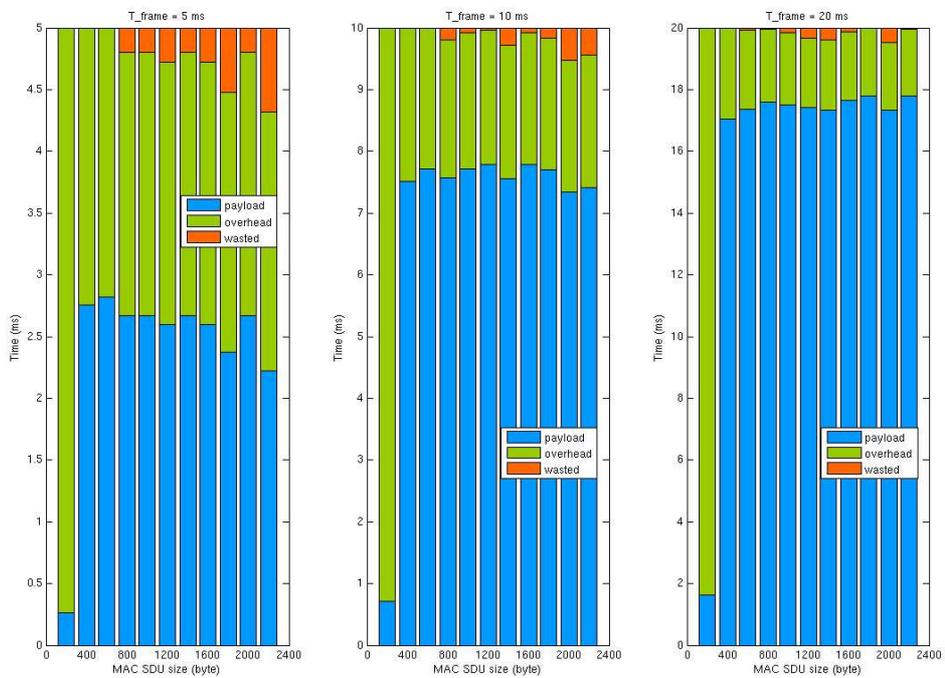


(b) 16-QAM 3/4

Figure 2.3: Effect of frame duration and modulation and coding scheme on IP throughput



(a) QPSK 1/2



(b) 64-QAM 3/4

Figure 2.4: Effect of the frame duration and the MCS on bandwidth utilization

### 2.2.2 Effect of the channel bandwidth

Recall that in previous scenarios, the channel bandwidth was fixed to 7 MHz. The scenarios considered in this section are aimed at showing the effect of the channel bandwidth on MAC goodput, therefore we will consider different values of channel bandwidth which implies different values of sampling factor (see Table 2.1) and consequently different durations of an OFDM symbol as we have seen in Section 2.1. However, we are more interested here in evaluating the MAC efficiency than in knowing the corresponding value of IP throughput. The MAC efficiency is defined as the percentage ratio between the MAC goodput (corresponding to the transmission of the MAC payload) and the physical rate.

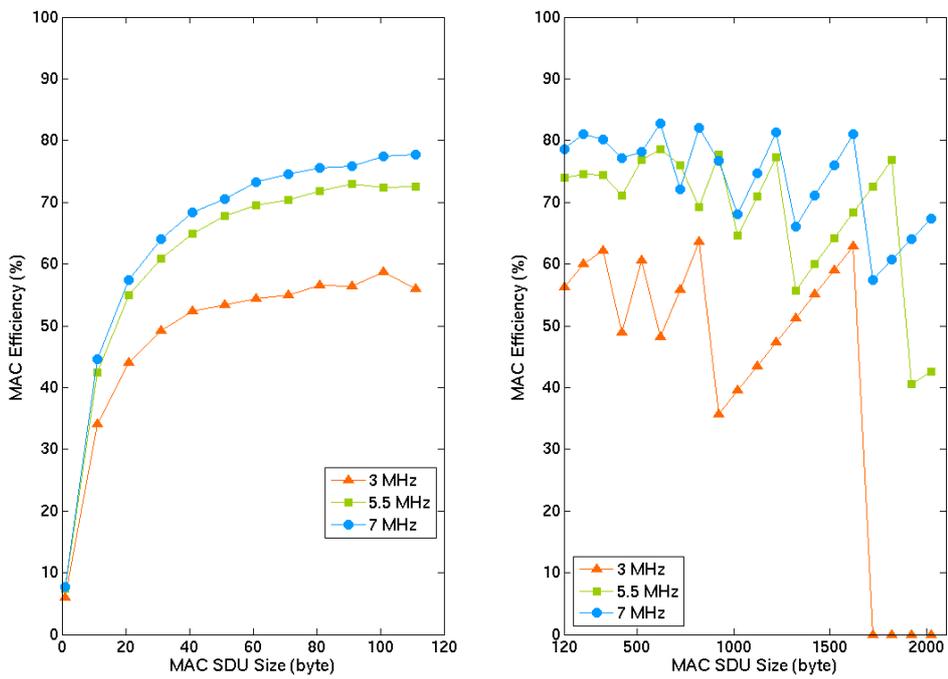
Figure 2.5(a) and 2.5(b) depict the MAC efficiency as a function of the MAC SDU size for three values of channel bandwidth: 3, 5.5, and 7 MHz. Note that each value corresponds to one of the PHY systems profiles specified by the IEEE 802.16 standard and reported in Table 2.1. The results presented in both figures are obtained for a frame duration of 10 ms. However in Figure 2.5(a), QPSK 1/2 is used while in 2.5(b) the modulation and coding scheme is set to 64-QAM 3/4.

Comparing the two figures, we observe that the obtained curves fluctuate a lot when using QPSK 1/2, and the larger is the bandwidth channel, the less visible is the fluctuation. This effect is similar to the one observed when varying the frame duration in 2.3(b) but here it is more discernible. In Figure 2.5(a), we see that for a channel bandwidth of 3 MHz, reaching a certain value of MAC SDU size (almost 1600 bytes), packet transmission is no longer possible with a frame duration of 10 ms. This is due not only to the shortness of channel bandwidth and frame duration but also to the absence of fragmentation. The two other curves corresponding to a channel size of 5.5 and 7 MHz, respectively exhibit almost the same behavior. Indeed, with MAC SDUs of more than 100 bytes, the MAC efficiency for 5.5 MHz fluctuates between 43.05 % and 80.84 % while for 7 MHz it varies between 56.5 % and 84.7 %. With a modulation and coding scheme of 64-QAM 3/4, the same behavior is observed since MAC efficiency fluctuates between 64.16 % and 73.29 % for a channel bandwidth of 5.5 MHz while it is between 71.66 % and 76.79 % for a channel bandwidth of 7MHz. The conclusion that may be derived from this is that the use of more than 20 % of extra bandwidth in the case of a channel size of 7 MHz does not imply a considerable improvement on MAC efficiency.

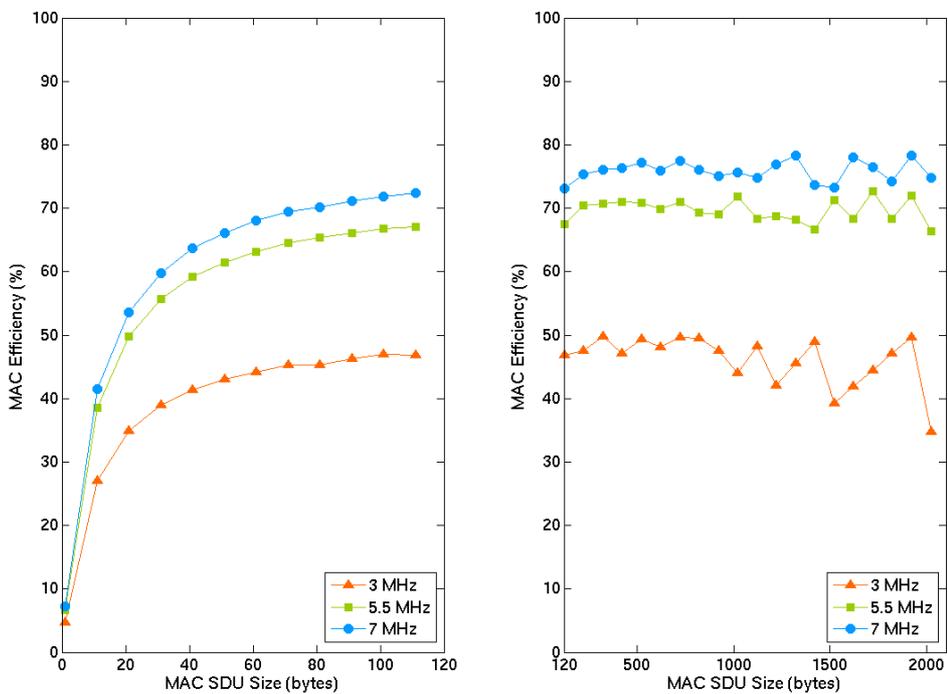
### 2.2.3 Impact of fragmentation and packing

Till now, the observed MAC performances were obtained when fragmentation and packing were disabled. However, we are interested in seeing how we could take advantage of these techniques, offered by the IEEE 802.16 standard [1], to improve the MAC efficiency. For this purpose, we consider the same plot shown in Figure 2.4(a) for a frame duration of 10 ms. Recall that this plot was obtained when both fragmentation and packing were deactivated. In the proposed scenario, we keep the same frame duration and MCS i.e. 10 ms and QPSK 1/2, respectively.

As fragmentation and packing are mutually exclusive [1], we first activate packing and prohibit fragmentation (see Figure 2.7). Note that we consider the fixed-length MAC SDUs variant of packing since the MAC SDUs have the same size. Comparing the proportions of overhead obtained when packing is activated and when not (Figure 2.7), we notice that packing has almost no impact on wasted ratio however it considerably increases the throughput when the MAC SDUs are small. This may be explained by the fact that when packing fixed-length MAC SDUs, only one packing subheader is needed for the whole MAC PDU what decreases considerably the resulting overhead particularly for small MAC SDUs (less than 400 bytes). Indeed instead of having a MAC header and a CRC field for each MAC SDU, we need only one generic MAC header, one CRC field, and a single packing subheader for all the MAC SDUs transmitted during a time frame.



(a) 10 ms and QPSK 1/2



(b) 10ms and 64-QAM 3/4

Figure 2.5: Effect of the channel bandwidth on MAC efficiency

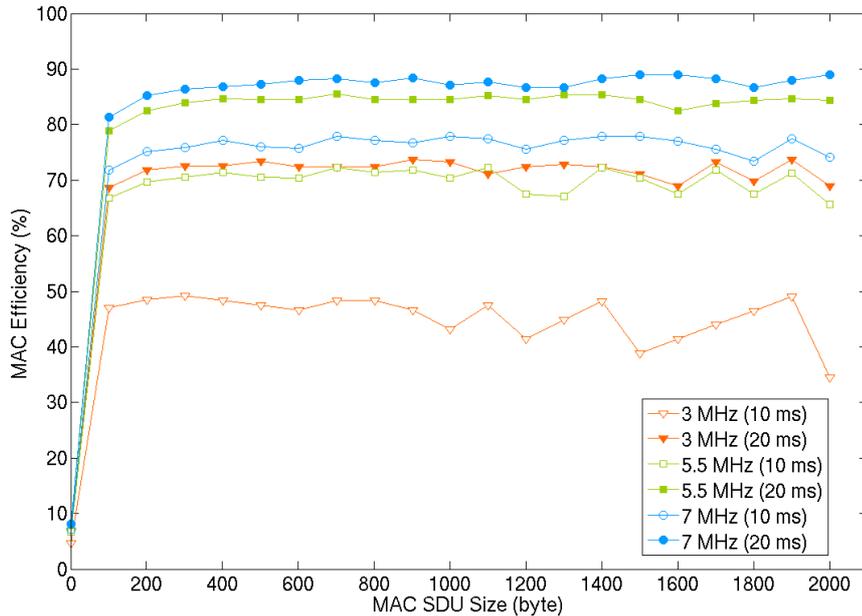


Figure 2.6: MAC efficiency using 64-QAM 3/4

Still referring to Figure 2.7, we are interested in seeing the impact of fragmentation on the frame composition. Comparing the case where packing and fragmentation are disabled to when the latter is enabled, we notice that the unused proportion of bandwidth is used to send more data and, of course, the resulting overhead. However the improvement of IP throughput is hardly discernible even though we are considering the optimal fragmentation case i.e. where the fragment size is adapted to the unused bandwidth. In Figure 2.6, we combine the variation of channel bandwidth along with the frame duration with 64-QAM 3/4 as modulation and coding scheme. It is interesting to note that all the curves have almost the same behavior. As expected, the larger is the channel and frame sizes, the higher is the MAC efficiency.

### 2.3 Conclusion

In this chapter, an original analytical framework was developed to investigate the performance bounds of OFDM-based 802.16 systems. This analytical framework was carried out with respect to what have been specified in the IEEE 802.16 standard [1]. It outlines a number of key features proposed by the standard and that have been hardly addressed in previous research works. Based on this framework, several scenarios were considered to evaluate the performance bounds of 802.16 systems under different MAC and PHY settings. The obtained results highlight the importance of considering the MAC and PHY overhead when evaluating the performance of IEEE 802.16 systems. Indeed this overhead, that is usually ignored or roughly estimated in most research works related to WiMAX resource allocation, may constitute 80 % of the whole frame. Also we have shown that using a larger bandwidth channel may yield minimal improvements on MAC performances. Also when investigating fragmentation and packing impact on MAC performance, we have shown that packing may considerably improve the resulting throughput especially for traffic carrying fixed-size packets.

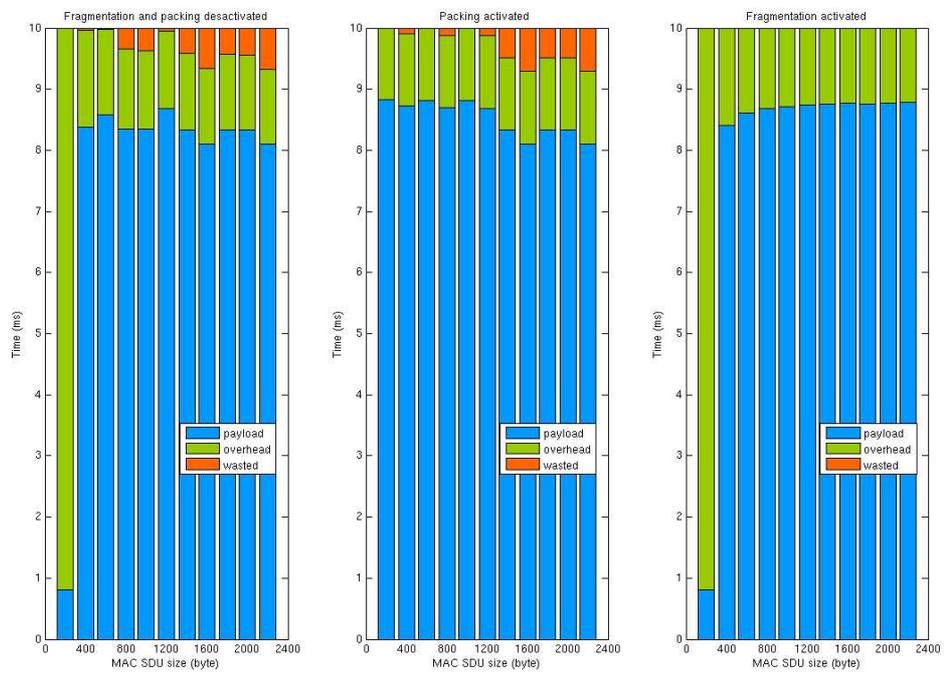


Figure 2.7: Effect of packing and fragmentation



## Chapter 3

# QoS Support in WiMAX Networks

The IEEE 802.16 standard defines a connection-oriented MAC protocol that is designed to accommodate a variety of applications with different QoS requirements. Nevertheless, several issues mainly related to resource allocation, have been left open. The main objective of this chapter is to provide a better understanding of the supported and missing features to ensure QoS support in IEEE 802.16 networks. Therefore, we describe in Section 3.1 the main elements specified by the IEEE 802.16 standard to provide QoS for heterogeneous classes of traffic. In Section 3.2, we propose a QoS architecture for WiMAX networks. The proposed architecture is intended to be a generic framework which incorporates what we consider as key components to answer the QoS needs of the different categories of applications addressed by the WiMAX technology. Section 3.3 is dedicated to scheduling and admission control issues. More specifically, we point out, through this section, the main challenges faced when designing a scheduling and/or CAC solution for WiMAX networks. Section 3.4 concludes the chapter.

### 3.1 QoS support in WiMAX networks

#### 3.1.1 Service flows management and QoS requirements

The standard defines a connection-oriented MAC protocol where all the transmissions occur within the context of a unidirectional connection. Each connection, identified by a unique Connection ID (CID), is associated to an admitted or active service flow (SF) whose characteristics provide the QoS requirements to apply for the protocol data units (PDUs) exchanged on that connection. In order to facilitate the MAC service data units (SDUs) delivery with the appropriate QoS constraints, the IEEE 802.16 Standard defines a classification process by which a MAC SDU is mapped to the associated connection and so to the SF corresponding to that connection. The classification procedure is performed at the service-specific convergence sublayer (CS) by classifiers consisting of a set of protocol-specific matching criteria (c.f. Section 1.4.2).

There are three types of service flows: (a) provisioned service flows for which the QoS parameters are provisioned for example by the network management system, (b) admitted service flows for which resources—mainly bandwidth—are reserved and (c) active service flows which are activated to carry traffic using resources actually provided. Each service flow is uniquely identified by a service flow identifier (SFID). Service flows may be dynamically managed. They may be created, changed or deleted using Dynamic Service Addition (DSA), DS Change (DSC), and DS Delete (DSD) MAC management messages, respectively. These operations could be initiated either by the BS (mandatory capability) or by the SS (optional capability). Figures 3.1(a) and 3.1(b) illustrate the two cases for the creation of a SF. Within these three/four-ways hand-

---

shakes, a Dynamic Service Addition Request (DSA-REQ), a DSA Response (DSA-RSP), and a DSA Acknowledgement (DSA-ACK) messages are exchanged between the BS and the SS. When the transaction (addition of a SF) is initiated by the SS, the BS transmits an extra message DSA Received (DSA-RVD) informing the SS that the DSA-REQ has been received and is being treated by the BS. The DSA-REQ includes:

- a Transaction ID, assigned by the sender, that uniquely identifies the current transaction.
- a set of service flow parameters specifying the flow's traffic characteristics and scheduling requirements.
- a SFID if the SF creation is initiated by the BS. In this case, the DSA-REQ may also include a CID when the SF is admitted.

The DSA-RSP transmitted in response to a DSA-REQ indicates the acceptance or rejection of the SF. A specific parameter called confirmation code (CC) specifies whether a SF was accepted or not and the cause of the rejection (e.g. CC = 0 indicates an OK/success, CC = 3 indicates the absence of sufficient resources to admit the SF [37]). If the DSA-RSP includes a newly assigned CID, it should also contain the complete set of QoS parameters. This set specifies for instance:

- the minimum reserved traffic rate: expressed in bits per second, this parameter indicates the minimum rate reserved for the service flow. When omitted, a default value of 0 is considered.
- the maximum sustained traffic rate: expressed in bits per second, it defines the peak information rate for the service.
- the maximum latency: it defines the maximum interval between the entry of a packet at the CS of the BS or the SS and the forwarding of the SDU to its Air Interface. If specified, the BS or SS is committed to guarantee it. Nevertheless, a BS or SS does not have to meet this service commitment for service flows that exceed their minimum reserved rate.
- the SDU size parameter: this parameter specifies the length of the SDU for a fixed-length SDU SF.

Instead of explicitly specifying the whole set of QoS parameters characterizing the SF, a flow can be created by specifying a service class name that identifies a set of QoS traffic parameters. The concept of a service class is an optional capability and may be implemented at the BS. It allows higher layers to instantiate a service only by specifying its service class name. For example, telephony signaling may direct the SS to instantiate any available provisioned service flow of class "G711" [37].

### **3.1.2 Scheduling service types**

Depending on the service to be tailored to each user application, a specific scheduling service is attributed to handle the flow. Based on that, a specific set of QoS parameters should be specified when creating a new service flow (like it is shown in Table 3.1). Uplink flows however are associated, in addition to a scheduling service, to one of these request/grant scheduling types: unsolicited grant service (UGS), real-time polling service (rtPS), extended real-time polling service (ertPS)—introduced by the IEEE 802.16e-2005 standard [2], non-real-time polling service (nrtPS), and best effort (BE). Each scheduling service is designed to meet the QoS requirements of a specific applications category. More details about each request/grant scheduling type are given in the next paragraphs.

---

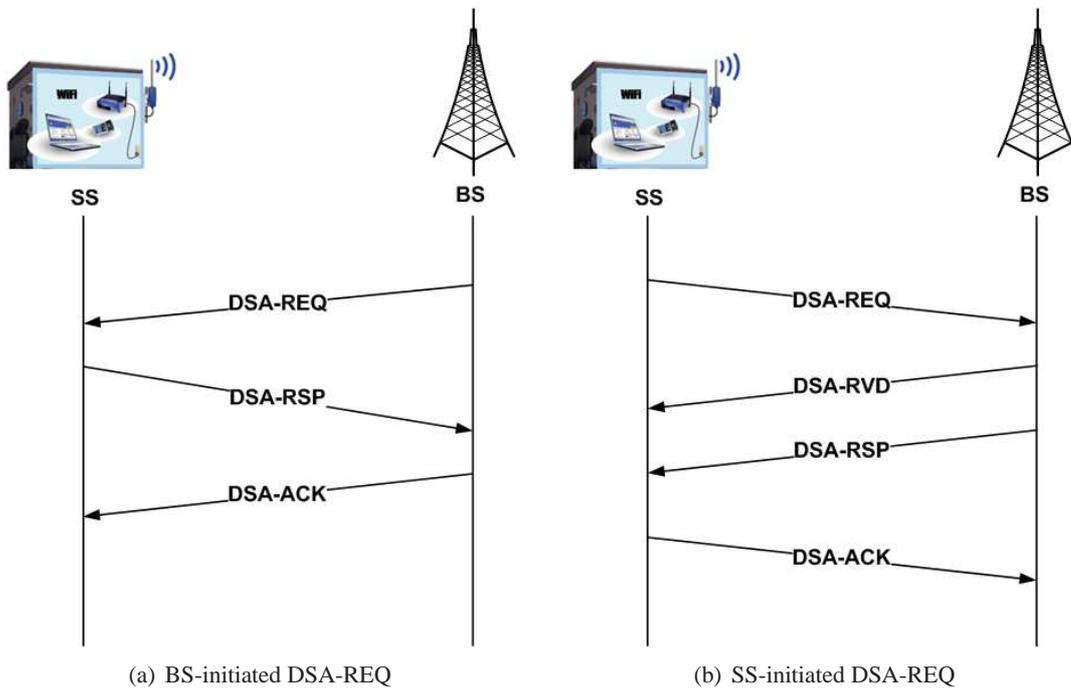


Figure 3.1: Dynamic Service Addition

- **UGS** is designed to support real-time applications that generate fixed-size data packets at periodic intervals, such as T1/E1 and voice over IP (VoIP) without voice activity detection (VAD). The mandatory service flow QoS parameters for UGS service are listed in Table 3.1. This table summarizes, according to the scheduling service type, the QoS parameters that must be specified when establishing a new service flow. UGS connections never request bandwidth. The amount of bandwidth to allocate to such connections is computed by the BS based on the minimum reserved traffic rate defined in the service flow of that connection.
- **rtPS** is designed to support real-time applications that generate variable-size data packets at periodic intervals, such as moving pictures expert group (MPEG) video. Unlike UGS connections, rtPS connections must inform the BS of their bandwidth requirements. Therefore the BS must periodically allocate bandwidth for rtPS connections specifically for the purpose of requesting bandwidth. This corresponds to the polling bandwidth-request mechanism. This mechanism exists in three variants: unicast polling, multicast polling and broadcast polling. Only unicast polling can be used for rtPS connections.
- **Extended rtPS** is a new scheduling service introduced by the IEEE 802.16e-2005 standard [2] to support real-time service flows that generate variable size data packets on a periodic basis, such as Voice over IP services with silence suppression. Like in UGS, the BS shall provide unicast grants in an unsolicited manner which saves the latency of a bandwidth request. However, unlike UGS allocations that are fixed in size, ertPS allocations are dynamic like in rtPS. By default, the size of allocations corresponds to the current value of Maximum Sustained Traffic Rate at the connection. The SS however may request changing the size of the UL allocation.
- **nrtPS** is designed to support delay-tolerant applications such as FTP for which a minimum amount of bandwidth is required. The polling mechanism can be applied to nrtPS connec-

Traffic/Applications Characteristics	real-time, fixed-rate data, Fixed/Variable length PDUs		real-time, variable bit rates, requiring guaranteed data rate and delay		real-time, variable bit rates, requiring guaranteed data rate and delay		requiring guaranteed data rate, insensitive to delays		No rate or delay requirement	
	DL	UL	DL	UL	DL	UL	DL	UL	DL	UL
Downlink (DL)/ Uplink (UL)										
Maximum Sustained Traffic Rate	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Minimum Reserved Traffic Rate	✓	✓	✓	✓	✓	✓	✓	✓	—	—
Maximum Latency	✓	✓	✓	✓	✓	✓	—	—	—	—
Tolerated Jitter	✓	✓	✓	✓	—	—	—	—	—	—
Request/Transmission Policy	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Traffic Priority	—	—	✓	✓	✓	✓	✓	✓	—	—
Request/Grant Scheduling Type	—	✓ (UGS)	—	✓ (ertPS)	—	✓ (rtPS)	—	✓ (nrtPS)	—	✓ (BE)
Unsolicited Grant Interval	—	✓	—	✓	—	—	—	—	—	—
Unsolicited Polling Interval	—	—	—	—	—	✓	—	—	—	—
SDU Size (If fixed length SDU)	✓	✓	—	—	—	—	—	—	—	—
Example of application	T1/E1, VoIP without VAD		VoIP with VAD		MPEG video		FTP		HTTP, SMTP	

Table 3.1: Mandatory QoS parameters for each scheduling service

tions. However, unlike for rtPS, nrtPS connections are not necessarily polled individually—multicast and broadcast polling are possible—and the polling must be regular, not necessarily periodic.

- **BE** is designed for applications that do not have any specific bandwidth or delay requirement, such as HTTP and SMTP. For BE connections, all forms of polling are allowed in order to request bandwidth.

The QoS parameters that must be specified when establishing a new service flow are listed in Table 3.1. The value of the Request/Transmission (Rx/Tx) Policy parameter offers the possibility to specify options for PDU formation. It might define for instance a restriction on packing and fragmentation capabilities as well as attributes affecting the bandwidth request types.

### 3.1.3 Bandwidth allocation and request mechanisms

Except for UGS connections that receive the bandwidth in an unsolicited manner, the MS needs to inform the BS of its uplink requirements. To do so, a set of mechanisms is proposed by the IEEE 802.16 standard.

- **Polling** To poll an SS, the BS allocates enough bandwidth to send a bandwidth request (BR). This bandwidth request opportunity is specified in the UL-MAP through a request information element (IE). There exist several forms of polling: unicast polling (addressed to the basic CID of an SS), multicast polling, group polling and broadcast polling. Nevertheless,

the use of the one or the other of these forms of polling is restricted by the scheduling service type of the considered connection. For ertPS and rtPS, only unicast polling is allowed. For nrtPS and BE, all forms of polling are possible.

- **Piggybacking** To ask for bandwidth, the SS may send a stand-alone bandwidth request header (6 bytes) or just piggyback the request on a PDU using a grant management (GM) subheader (2 bytes). The support of piggybacking is optional and may be used only to request bandwidth for the connection carrying the PDU to which the GM subheader has been added.
- **Bandwidth stealing** This mechanism refers to the use, by the MS, of a portion of the bandwidth allocated for data (through a data grant IE) to transmit a bandwidth request instead.
- **PM-bit** MSs having at least one active UGS connection may set the poll-me (PM) bit, of the GM subheader, in a MAC PDU of the UGS connection to inform the BS that polling is needed for non-UGS connection. As a response to this poll request, the BS initiates a process of unicast polling. As specified in [37], this technique should be used by the MS only when piggybacking cannot be performed and if all the possibilities of bandwidth stealing are exhausted.

It is worth mentioning that, whatever is the bandwidth request mechanism in use, bandwidth is always requested by an SS on a per-connection basis and addressed by the BS to the SS as an aggregate of grants. Therefore, since the SS receives the allocated bandwidth as a whole in response to per-connection requests, it cannot know which request is honored. The SS can then use the grant either to send data, or to request bandwidth for any of its connections (bandwidth stealing), or even to send management messages.

## 3.2 A QoS architecture for WiMAX networks: the big picture

The framework we propose in this section is independent of the adopted scheduling and CAC strategy. It is the compilation of what we consider as key elements for QoS support in WiMAX systems. The names, roles and interactions between the different entities described in this section represent a proposal among others for a QoS framework addressing WiMAX systems. Figure 3.2 illustrates the different elements and modules that constitute the proposed framework as well as the interactions that exist between them. In this figure, we can see the proposed MAC logical structures for both BS and SS. Note that entities having the same name appear in both sides. In general, they play the same role. Differences will nevertheless be shown and discussed while explaining the role of each component.

- *Classifier*: As mentioned in Sections 1.4.2 and 3.1, when a MAC SDU is received, it should be mapped to a particular connection. In the proposed framework, this task is accomplished by the *Classifier* based on a set of matching criteria such as the 5-IPv4 tuple<sup>1</sup>. As can be seen from Figure 3.2, the classification process is applied by both BS and SS to packets they are transmitting. When a BS or an SS receives an SDU, it refers to a matching table; if the considered SDU matches the criteria relative to a specific CID, it is then transmitted to the *Buffer Manager* (see Figure 3.2).

<sup>1</sup>IP source and destination addresses, source and destination ports, and QoS type field

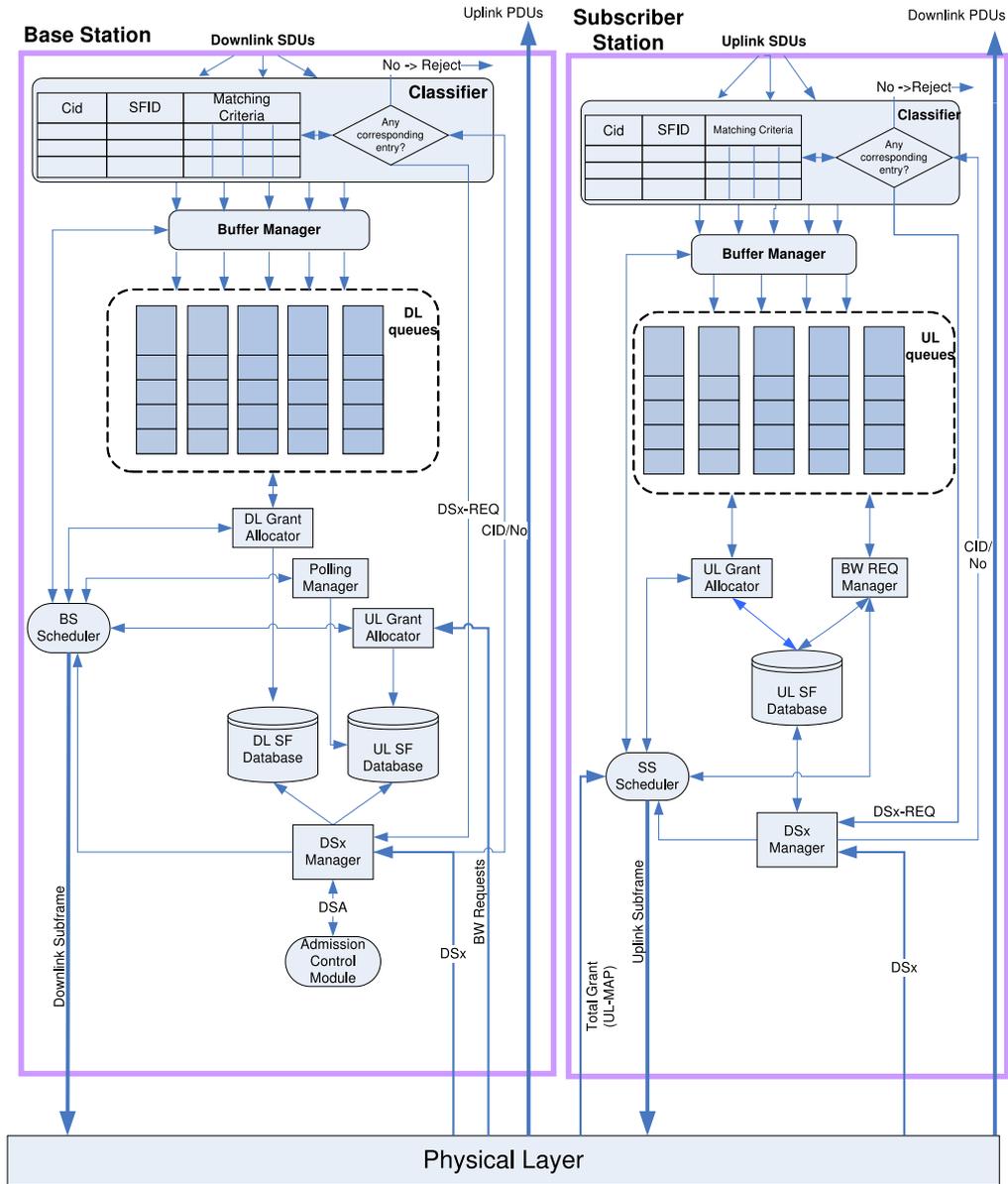


Figure 3.2: QoS architecture Design

It is possible that the SDU cannot be associated to any active connection. The *Classifier* checks then whether the packet can be mapped to a provisioned service flow (a SF that does not correspond to any CID), in which case the *Classifier* asks the *DSx Manager* to activate the corresponding SF. The *DSx Manager* may either respond by specifying a new CID or reject the Classifier demand—generally as a result of *Admission Control Module* decision (see Figure 3.2). In the latter case, as well as when the MAC SDU cannot be associated to any matching table entry, the packet is dropped.

- *Buffer Manager*: The *Buffer Manager* is responsible for managing the queues. It allocates a separate queue for each active MAC connection. The *Buffer Manager* operates as follows. When it receives a MAC SDU from the *Classifier*, it puts it into the corresponding connection queue based on the CID already determined by the *Classifier*. It may also discard the packet if the buffer capacity of the corresponding connection is exceeded. As we will see later, the *Buffer Manager* is also responsible for delivering a specific amount of data from each connection queue according to what has been decided during the scheduling procedure.
- *DL/UL SF Database*: This database contains the QoS parameters of each DL/UL service flow. These parameters depend on the service flow type: whether it is associated to a UGS, ertPS, rtPS, nrtPS or BE connection. The *DL* and *UL SF Databases* are used, during the scheduling procedure, to get the DL and UL QoS constraints, respectively. These databases are managed and maintained up-to-date by the *DSx Manager*.
- *DSx Manager*: The *DSx Manager* receives and treats all the messages exchanged during a service flow creation, deletion or change procedure. For example, when the BS receives a DSD-REQ message, the *DSx Manager* treats the request and informs the *BS Scheduler* that a DSD-RVD message should be sent during the next frame interval. DSA-REQ messages, and in some cases DSC-REQ messages should be handled differently. In fact, since the purpose behind sending such messages is to create a new active service flow (DSA-REQ) or to activate an existing provisioned service flow (DSC-REQ), the request shall be handled by the *Admission Control Module* (see Figure 3.2) which decides whether it can be accepted or not. According to the request and to the *Admission Control Module* decision—when considering the case of a service flow creation or activation—the *DSx Manager* updates the *DL/UL SF Database* by entering, modifying or deleting the QoS characteristics of the considered service flow.
- *Admission Control Module*: Note that, although it exists only at the BS side (see Figure 3.2), the *Admission Control Module* is applied for both SS and BS-initiated connections addition requests. Its role is to check whether a request to create a new active service flow or to activate an existing one can be honored; other cases like using a more robust modulation technique may also require *Admission Control Module* action. The decisions made by the *Admission Control Module* should be communicated to the *DSx Manager* which is responsible for planning and applying them.
- *Polling Manager*: As mentioned in Section 3.1.3, [37] proposes several bandwidth request mechanisms. Apart from UGS which does not need any explicit bandwidth request, polling can be applied to all types of scheduling services. Therefore, we integrate in the proposed architecture design the *Polling Manager* (see Figure 3.2) whose role consists in granting bandwidth request opportunities. The *Polling Manager* needs to have access to the *UL SF Database* to determine the SSs that are concerned by this technique, which

are those having at least one non-UGS connection. As specified in the IEEE 802.16 standard [37], polling should be periodic for rtPS connections and regular for nrtPS connections.

- *DL Grant Allocator*: The *DL Grant Allocator* allocates bandwidth for DL connections according to the adopted scheduling policy. To accomplish this task, the *DL Grant Allocator* needs to know the current status of DL queues, the QoS constraints of each DL service flow and the amount of remaining bandwidth—a parameter that is maintained by the *BS Scheduler*. To get the DL SF parameters, the *DL Grant Allocator* refers to the *DL SF Database*.
- *UL Grant Allocator*: As shown in Figure 3.2, this entity exists in both BS and SS structures.
  - At the BS, bandwidth is allocated based on bandwidth requests sent by the SSs for polling services and BE connections. For UGS connections, the grants are made according to the QoS constraints of the associated service flow. In all cases, the *UL Grant Allocator* needs to have access to the *UL SF Database* to get the QoS constraints of UL service flows. After allocating bandwidth, the *UL Grant Allocator* informs the *BS Scheduler* of its grants decisions.
  - At the SS, the *UL Grant Allocator* operates in the same manner as the *DL Grant Allocator* at the BS with the exception of dealing with UL connections instead of DL connections.
- *BW REQ Manager*: As mentioned above, the *Polling Manager* has to allocate bandwidth to an SS specifically for the purpose of requesting bandwidth for its non-UGS UL connections. These allocations are then used by the *BW REQ Manager*, at the SS, to send bandwidth requests; they may optionally be used to send data. More generally, the SS may use any uplink allocation to send data or bandwidth requests [37]. Therefore, since the *BW REQ Manager* does not know the exact amount of bandwidth to be used for requests, it should refer to the *SS Scheduler* which is the only component able to make such decisions and to have information on bandwidth availability. Also, the *BW REQ Manager* needs to check the UL queues and *UL SF Database* in order to plan the requests.
- *BS Scheduler*: The *BS Scheduler* represents the main element of the proposed architecture. In fact, it is responsible for coordinating the work of the *DL Grant Allocator*, the *UL Grant Allocator* and the *Polling Manager* since it maintains information on the amount of the remaining bandwidth after each scheduling step. Besides, the *BS Scheduler* should remain informed of the *DSx Manager* decisions in order to plan the DSx management messages, such as DSA-REQ, DSC-RVD, DSD-RSP, to be sent in the current frame. Based on all the collected information, the *BS Scheduler* first generates the DLFP, the UL-MAP, and optionally the DL-MAP messages. These messages hold the scheduling decisions made by the BS (more specifically by the *DL Grant Allocator*, the *UL Grant Allocator* and the *Polling Manager*). Secondly, it either asks the *Buffer Manager* to transmit data according to what has been specified in the DLFP message, or just generates the appropriate management messages and send them to the SSs.
- *SS Scheduler*: As far as the SS is concerned, the *SS Scheduler* is the main element in the proposed QoS architecture. The *SS Scheduler* interacts with the BW-REQ

Manager and the *UL Grant Allocator* in order to use the whole grant assigned by the BS. Note that only the *SS Scheduler* knows the amount of bandwidth that was granted to the SS's Basic CID. Thus, it has the responsibility of updating this parameter after each scheduling decision.

In general, any QoS framework addressing WiMAX systems should consist of three main building blocks: a DL and UL scheduler at the BS, an UL scheduler at the SS and possibly an admission control module. Nevertheless, the details of these blocks are vendor-specific.

### 3.3 Scheduling and CAC in WiMAX: design challenges

The objective of this section is to provide a better understanding of the design challenges of a new scheduling and/or CAC solution for IEEE 802.16 since they represent the main issues for insuring QoS.

- **QoS requirements guarantee:** The scheduler should satisfy the QoS requirements of the different types of service specified by the standard. Hence it has to monitor, for each connection, the required QoS parameters, presented in Table 3.1, and check if they are in line with what has been negotiated.
- **Bandwidth-request strategy:** Because the standard gives a choice among several bandwidth request and grant techniques, it is important for each scheduling solution to define its own bandwidth request strategy.
- **Graceful service degradation:** It is an interesting characteristic for CAC and scheduling algorithms, when accepting new connections, to degrade the service of the ongoing over provisioned connections as gracefully as possible. Since radio resources are limited the use of this kind of strategy would compensate lagging flows and ensure fairness in radio resources management (RRM).
- **Channel utilization:** The channel utilization is expressed in percentage of the available capacity and it represents the achieved throughput. It corresponds to the fraction of time used to transmit data packets. In the case of a PMP communication, this parameter is almost equal to the channel capacity. Nevertheless, to maximize the channel utilization, the scheduler should minimize the overhead by optimizing the bandwidth-request strategy and taking advantage of the concatenation, packing, and fragmentation mechanisms, proposed by the standard.
- **MAC-PHY cross-layer design:** This constraint consists mainly in considering the adaptive modulation and coding (AMC) capability defined by the standard. Indeed, it is important, when allocating resources at the MAC level, to take into account the burst profile in use.
- **Fairness:** One of the most challenging problems for RRM is to find a compromise between increasing the channel utilization— by serving flows with good channel conditions— and being fair to different flows. To estimate this parameter Jain's fairness index might be used:

$$F_J = \frac{(\sum_{i=1}^m x_i)^2}{m \cdot \sum_{i=1}^m x_i^2}$$

Where  $m$  is the total number of flows and  $x_i$  is the proportion of received packets of flow  $i$  during run time.  $F_J$  is equal to 1 when all flows equally share the bandwidth, and equal to  $1/m$  when a flow monopolizes the network.

- **Implementation complexity:** Scheduling and CAC algorithms deal with many different constraints. Nevertheless, because they address—among others—real time flows, they need to be fast and should not have a prohibitive implementation complexity.
- **Scalability:** Scalability is the capability of the scheduling algorithm to handle growing number of flows, or nodes, in a graceful manner. Scalability is also important in the context of mobile WiMAX networks for mobility management.

### 3.4 Conclusion

From this chapter we have seen that the IEEE 802.16 standard defines:

1. Concepts making easier to associate packets with the appropriate QoS constraints, namely concepts of connections, service flows, classes of services, and classifiers;
2. Five scheduling services tailored to meet the QoS requirements of heterogeneous classes of traffic;
3. Signaling mechanisms offering the possibility to manage service flows dynamically (i.e. DSx messages) and to request (e.g. piggybacking) and grant (e.g. UL-MAP message) bandwidth;
4. A scheduling procedure for UGS connections.

The standard, nevertheless, leaves undefined:

1. The admission control policy to apply when the creation of a new service flow (or the activation of a provisioned one) is requested.
2. The scheduling mechanisms based on which resources shall be allocated.

Based on the above considerations, we have proposed in this chapter a generic QoS framework which incorporates the main supported and missing functionalities to handle QoS in WiMAX systems. We have tried, when designing this framework to be as close as possible to what has been specified by the IEEE 802.16 standard [37].

The last section of this chapter has been dedicated to scheduling and admission control issues. More specifically, we have highlighted, through that section, the main challenges faced when designing a scheduling and/or CAC solution for WiMAX networks. These constraints represent also the main evaluation criteria of the different resource management mechanisms proposed in this work-in progress area. The state of the art of these mechanisms is presented in next chapter.

---

## Chapter 4

# Scheduling and CAC in WiMAX Networks: a Survey and Taxonomy

A large body of literature has been concerned with scheduling and admission control issues in WiMAX networks. In this chapter, we survey, classify, and compare different scheduling and CAC mechanisms proposed in this work-in-progress area. The remainder of this chapter is divided into two main sections: Section 4.1 and Section 4.2 provide a survey and taxonomy of scheduling and CAC mechanisms dedicated to WiMAX networks. Section 4.3 concludes the chapter by outlining the main concerns worth addressing in this field. Most of the works presented in this chapter are proposed for OFDM-based WiMAX networks.

### 4.1 Scheduling

As shown in Figure 4.1, the approaches adopted in literature when designing a scheduling solution can be divided into three main categories. (1) The first one is a queuing-derived strategy where the authors focus on the queuing aspect of the scheduling problem and try to find the appropriate queuing discipline that meet the QoS requirements of the service classes supported by the IEEE 802.16 standard [1, 2]. In this first category, two kinds of structures are proposed: either simple structures consisting in general in one queuing discipline applied for all the scheduling services [3, 4, 5] or hierarchical structures consisting in two or multiple layers reflecting different levels of scheduling like in [6, 7, 8, 9, 10, 11, 12, 13, 14]. (2) In the second category, the scheduling problem is formulated as an optimization problem whose objective is to maximize the system performance subject to constraints reflecting in general the QoS requirements of different service classes [15, 16, 17, 18, 19, 20, 21, 22, 23]. (3) The third category of scheduling mechanisms that can be found in literature is the cross-layer strategy. The scheduling schemes adopting this strategy are usually based on a cross-layer architecture. The objective of this architecture is to optimize the communication between two [24, 25, 26, 27, 28] or three different layers [29, 30] and thus improve the system performance. As we will see in Section 4.1.3, these schemes could be further classified based on the layers involved in the cross-layer design.

#### 4.1.1 Packet queuing-derived strategies

##### 4.1.1.1 One-layer scheduling structures

Sayenko *et al* [5] consider that because there is not much time to do the scheduling decision, a simple one-level scheduling mechanism is much better than a hierarchical one. Therefore they

---

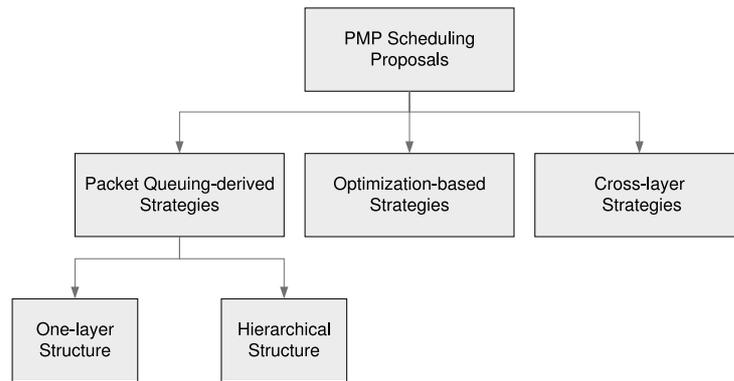


Figure 4.1: Classification of the scheduling strategies of IEEE 802.16 PMP mode

propose a scheduling solution based on the round-robin (RR) approach. They argue that there is no need to use disciplines like fair queuing (FQ) since the weights in such algorithms are floating numbers while the number of allocated slots, in 802.16 networks, should have an integer value. They also try to outline the difference between the weighted round-robin (WRR) discipline and the 802.16 environment. They insist on the fact that WRR may lead to a waste of resources because of its work-conserving behavior that does not fit the fixed-size frame of 802.16 that implies a non-work conserving behavior.

Based on the above considerations, the authors proposed in [5] a scheduling solution that consists in four main steps:

- Allocating for each connection the minimum number of slots that ensure the minimum reserved traffic rate with respect to the used modulation and coding scheme,
- Distributing the free slots between rtPS and nrtPS connections and then assigning the remaining to BE connections,
- Ordering the slots in such a manner the delay and jitter values are decreased.
- Estimating the overhead for UGS, ertPS, and in some cases nrtPS connections. This is not possible for rtPS and BE connections since it is more likely that the SDU size varies.

Note that [5] is one of the rare research works in which the overhead resulting from the scheduling decision, and packing or fragmentation capability is taken into account. However it is also worth mentioning that the authors consider a grant per connection (GPC<sup>1</sup>) mechanism and when ordering slots, they apply an interleaved scheme that is in contradiction with the frame structure specified by the standard.

In [3, 4], Cicconetti *et al* conjecture that the class of latency-rate ( $LR$ ) scheduling algorithms is particularly suited for implementing schedulers in 802.16 MAC since the basic QoS parameter required by a given connection is the minimum reserved traffic rate. Indeed the behavior of such algorithms is determined by two parameters which are the latency and the allocated rate [41].

<sup>1</sup>This approach consists in allocating the bandwidth on a per connection basis. In contrast with GPC, the grant per subscriber station (GPSS) refers to the allocation of bandwidth per SS. Both concepts should have been disused since the publication of the IEEE 802.16a-2003 Standard [40]. Indeed, it is clearly specified in [1, 2] that bandwidth is requested on a per connection basis while grants are aggregated and addressed as a whole for each SS.

From this class, the authors have chosen the deficit round robin (DRR) algorithm. DRR is simple to implement ( $O(1)$  complexity if specific allocation constraints are met) and provides, according to [3, 4], fair queuing in presence of variable length packets<sup>2</sup>. It nevertheless requires a minimum rate to be reserved for each packet flow; so even BE connections should be guaranteed a minimum rate. Also since this algorithm assumes that the size of the head-of-line packet is known, it can not be applied by the BS to schedule uplink transmissions. For this reason the authors have made the choice of implementing it as SS scheduler and as a downlink scheduler at the BS, since both BS and SS know the head-of-line packet sizes of their respective queues. To schedule uplink transmissions at the BS—based on backlog estimation—they have selected the WRR algorithm which belongs, like DRR, to the class of  $LR$  algorithms.

The simulation study carried by Cicconetti *et al* [3] demonstrated that the performance of 802.16 systems, in terms of throughput and delay, depends on several metrics such as frame duration, the mechanisms used to request UL bandwidth, the offered load partitioning—how traffic is distributed among SSs, the connections within each SS, and the traffic sources within each connection.

#### 4.1.1.2 Hierarchical scheduling structures

Wongthavarawat *et al.* [13, 14] are the first authors who introduced a hierarchical structure of bandwidth allocation for 802.16 systems. This hierarchical scheduling structure, shown in Figure 4.2, combines strict priority policy, among the service classes, and an appropriate queuing management discipline for each class: earliest deadline first (EDF) for rtPS, and weighted fair queuing (WFQ) for nrtPS. Fixed time duration is allocated to UGS connections and remaining bandwidth is equally shared among BE connections. In order to avoid starvation for lower priority connections, a policing module is included in each SS. It forces each connection to respect the traffic contract when demanding bandwidth. The proposed scheduling algorithm takes into account the queue size information and the service actually received by each connection. It also considers the arrival time and the deadline requirements of rtPS connections. However, the authors focused only on UL scheduling. They considered TDD mode and assumed that the durations of UL and DL subframes are dynamically determined by the BS but they did not specify how these proportions are fixed. The QoS architecture they proposed in [13] includes a token-bucket based admission control module that will be described in Section 4.2.

Most of the works that we will present in this section are "quite similar" to the scheduling model introduced by Wongthavarawat *et al.* in [13, 14]. Nevertheless, since more or less features are supported by each scheme, we have grouped them based on their main common contribution.

**Delay-aware scheduling** In [12], Sun *et al.* proposed a two-layers scheduling structure composed of a BS scheduler and an SS scheduler. At BS scheduler, priority is given to schedule data grants for UGS connections and bandwidth request opportunities for rtPS and nrtPS connections. The amount of bandwidth allocated in this phase is reserved during connections setup. Data grants for rtPS, nrtPS are then scheduled taking into account the information contained into bandwidth request messages and their minimum requirements. Finally, the residual bandwidth, if any, is redistributed in proportion to pre-assigned connections weights. The proposed SS scheduler considers a fixed priority scheme—1, 2, 3 and 4 for BE, nrtPS, rtPS and UGS scheduling service, respectively. Bandwidth is firstly guaranteed for UGS connections. rtPS packets are then scheduled based on their respective deadline stamps—corresponding to their *arrival\_time + tolerated\_delay*. Each

---

<sup>2</sup>This is in contradiction to what has been stated by Fattah and Leung in [42] where they qualify the fairness of DRR algorithm as "poor".

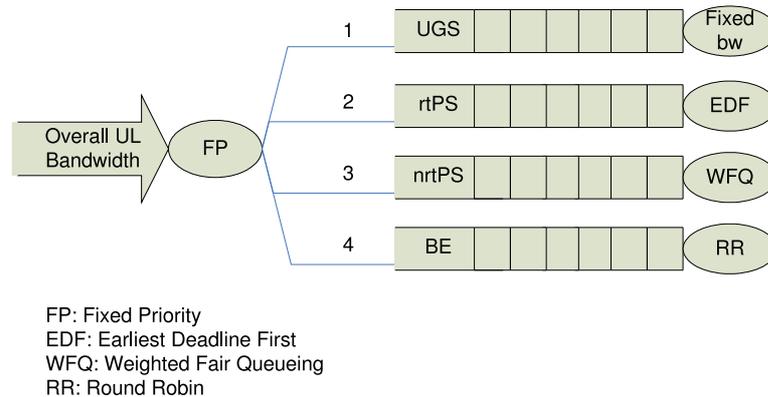


Figure 4.2: Hierarchical structure for bandwidth allocation [13, 14]

nrtPS packet is associated with a virtual time calculated to guarantee the minimum reserved bandwidth and hence maintain an acceptable throughput. A simple first-in-first-out (FIFO) mechanism is applied for BE queues.

Other scheduling schemes focusing on delay requirements were proposed in literature. In [8] for instance, three schedulers were combined to meet the QoS requirements of different classes (cf. Figure 4.3). Time sensitive traffic streams—namely UGS flows, rtPS flows and (n)rtPS polling flows—are served by Scheduler 1 that applies EDF algorithm. Minimum bandwidth reserving flows (nrtPS flows) are scheduled by Scheduler 2 using WFQ. The weights correspond to the proportion of requested bandwidth. WFQ algorithm is also applied by Scheduler 3 to serve BE traffics; weights nevertheless correspond in that case to traffic priorities specified by each BE connection. Other components of the proposed architecture are then used to plan contention and reserved transmission opportunities according to the bandwidth availability and to the priorities assigned to each scheduler—the highest priority is assigned to Scheduler 1.

In [10], a multimedia supported uplink scheduler is proposed by Perumalraja *et al.*. It includes a proportional fair (PF) BS scheduler and an earliest due date (EDD) SS scheduler. The BS scheduler (Figure 4.4.a) allocates resources first for the UGS service and then to poll SSs having at least one non-UGS connection: one slot is allocated in each frame for each SS having rtPS or nrtPS connections and one slot every three frames is allocated for SSs having only BE service connections. Finally, remaining OFDMA resources are proportionally allocated for SSs based on the received bandwidth requests. As can be seen from Figure 4.4.b, the EDD SS scheduler serves packets from the four traffic queues (UGS, rtPS, nrtPS and BE) in the order of the deadline assigned to each packet regardless of their scheduling service type.

**Asymmetric DL/UL scheduling** [7] is one of the rare research works that have proposed a scheduling algorithm considering simultaneously uplink and downlink bandwidth allocation in TDD mode. In first layer scheduling—of the two-layer hierarchical scheduling structure proposed in this work—Chen *et al* [7] have suggested the use of deficit fair priority queuing (DFPQ) algorithm instead of strict priority in order to avoid starvation for low priority classes. This first layer scheduling is based on two policies. The first one is a transmission direction-based priority where they chose to attribute to DL a higher priority than UL. The second policy is a service class-based priority applying the following scheme: rtPS>nrtPS>BE. As can be seen from Figure 4.6, the authors have combined these two policies using a strict priority scheme which assigns strict priority

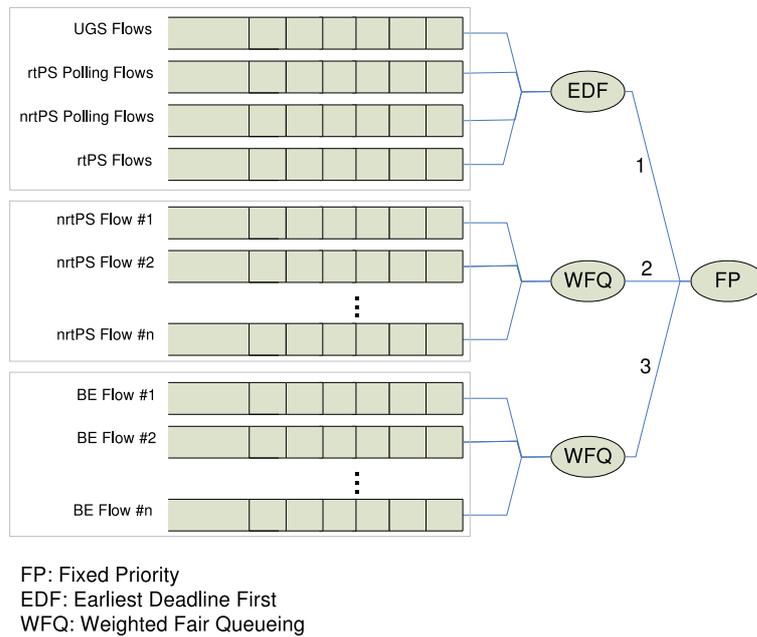
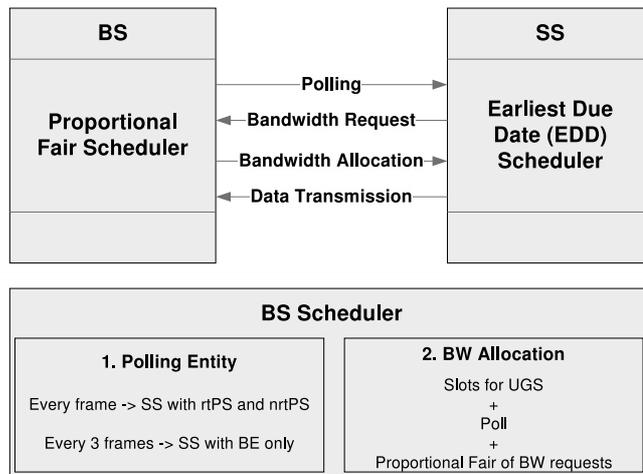
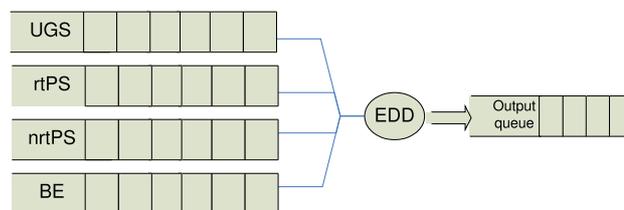


Figure 4.3: 3 schedulers proposal[8]



(a) BS scheduler [10]



EDD: Earliest Due Date

(b) EDD SS scheduler [10]

Figure 4.4: Multimedia supported uplink scheduler [10]

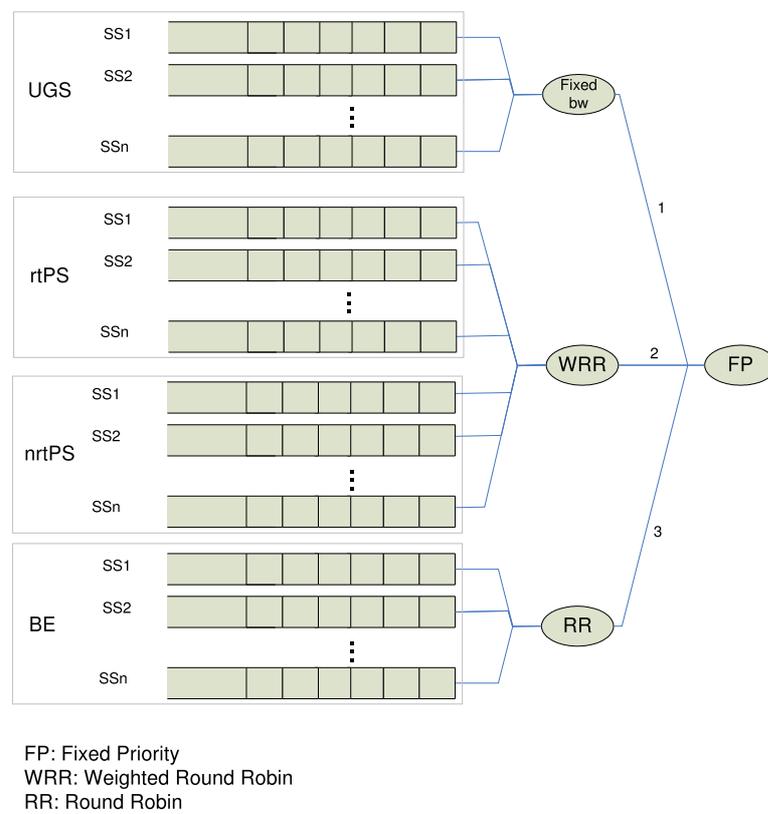
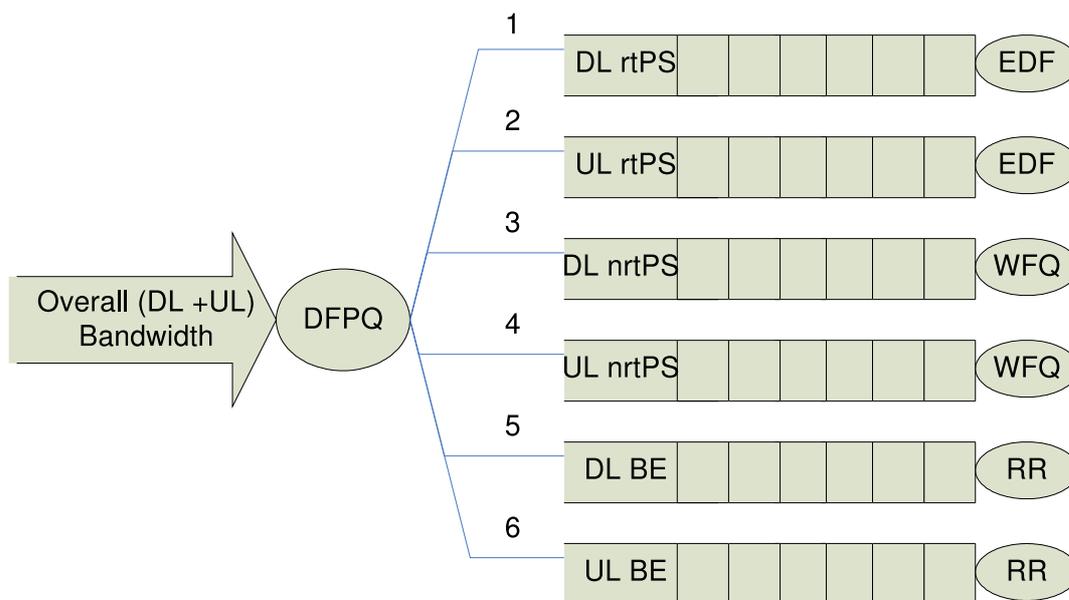


Figure 4.5: Scheduler model for WiMAX[11]

from highest to lowest to:  $DL_{rtPS}$ ,  $UL_{rtPS}$ ,  $DL_{nrtPS}$ ,  $UL_{nrtPS}$ ,  $DL_{BE}$ , and  $UL_{BE}$ . For DL and UL UGS connections, they have chosen to apply a fixed bandwidth allocation strategy. In second layer scheduling, three different algorithms were assigned to the other classes of services: EDF for rtPS, WFQ for nrtPS and RR for BE. nrtPS connections are scheduled based on weights corresponding to the ratio between the nrtPS connection minimum reserved traffic rate and the sum of the minimum reserved traffic rates of all nrtPS connections. A basic admission control algorithm is also proposed in this work. It accepts the connections for which the minimum reserved traffic rate does not exceed the available channel capacity; all BE connections are nevertheless accepted.

In order to take advantage of the DL/UL map of the 802.16d standard [1], Ma *et al.* propose in [9] a three-tier scheduling framework in which DL and UL respective loads could be unbalanced. Unlike in [7] however, the ratio of DL subframe with respect to the frame size is computed at the beginning of each frame. Indeed, a pre-scale dynamic resource reservation (PDRR) is used to allocate dynamically the overall frame bandwidth to DL and UL subframes with respect to a pre-scaled bound. The ratio of each subframe to the entire frame is computed based on the queues lengths and on the sizes of the bandwidth requests.



DFPQ: Deficit Fair Priority Queueing  
 EDF: Earliest Deadline First  
 WFQ: Weighted Fair Queueing  
 RR: Round Robin

Figure 4.6: Hierarchical structure of bandwidth allocation for WiMAX PMP mode [7]

**Packet-based scheduling: use of packing, fragmentation, PHS and AMC** Fragmentation, packing and PHS capabilities as well as their impact on the scheduling performance were considered in the packet-based scheduling strategy proposed in [11] by Settembre *et al.*. As can be

Scheduling proposal	Layer/Phase		DL	UL	UGS	rtPS	nrtPS	BE
[13, 14]	1 <sup>st</sup> layer				Fixed Priority			
	2 <sup>nd</sup> layer				Fixed Bandwidth	EDF	WFQ	Equally distributed
[12]	BS Scheduler	1 <sup>st</sup> phase		•	Fixed Bandwidth	Grant Bandwidth Request Opportunities		—
		2 <sup>nd</sup> phase			—	Guarantee the Minimum Reserved Rate		—
		3 <sup>rd</sup> phase			—	WFQ to distribute residual bandwidth		
	SS Scheduler				•	Fixed Priority		
					Fixed bandwidth	EDF	EDF (Virtual Time)	FIFO
[10]	BS Scheduler	1 <sup>st</sup> phase		•	Fixed Bandwidth	Unicast Polling		
		2 <sup>nd</sup> phase			—	Proportional Fair based on bandwidth Requests		
	SS Scheduler				EDD			
[11]	1 <sup>st</sup> layer		•		Fixed Priority			
	2 <sup>nd</sup> layer		•		Fixed Bandwidth	WRR		RR
[7]	1 <sup>st</sup> layer		•		DFPQ			
	2 <sup>nd</sup> layer		•		Fixed Bandwidth	EDF	WFQ	RR
[9]	Tier 1 (at BS)			•	Fixed Bandwidth	PQLW + MMFS among SSs		
	Tier 2 (at SS)			•	Fixed Bandwidth	SCFQ		WRR
	Tier 3 (per traffic flow)			—	EDF		SPLF	
[8]	Scheduler 1				EDF (UGS + rtPS + Polling rtPS and nrtPS)		—	—
	Scheduler 2				—	—	WFQ (based on bandwidth requests)	—
	Scheduler 3				—	—	—	WFQ (based on traffic priority)

Table 4.1: WiMAX hierarchical scheduling structures

seen from Figure 4.5, the proposed scheduler combines a strict priority policy among the different service categories and a specific queuing management discipline for each class: fixed bandwidth, WRR and RR for UGS, (n)rtPS and BE, respectively. For WRR discipline, weights are determined according to the guaranteed bandwidth.

Adaptive modulation and coding was also addressed in [11]. A preliminary WRR/RR allocation is achieved assuming the use of the most robust burst profile while bandwidth is allocated taking into account the actual burst profile! It is true that this way of proceeding guarantees enough bandwidth for existing flows even in the worst case. However, it might cause an unjustified high blocking rate and a low link utilization when the channel is good. Another shortcoming of [11] is that the admission control algorithm that manages the access of new connection—and based on which the minimum bandwidth requirements are guaranteed—is not described.

Table 4.1 summarizes the hierarchical scheduling proposals described above. In this table, we show whether DL connections are concerned or not by the proposed scheduling mechanism. Also, the table reflects the different steps of each scheduling process as well as the queuing discipline applied at each considered level of aggregation (per service type, per connection, etc.).

**Satisfaction-based scheduling** In [6], an original two-tier scheduling algorithm (2TSA) was proposed to avoid starvation problem and to provide fair allocation of residual bandwidth. UGS connection is not concerned by the “2TSA” algorithm since it is allocated a fixed amount of bandwidth per frame. Each connection is classified into either “unsatisfied”, “satisfied”, or “over-satisfied” connection and is assigned a weight indicating its shortage or satisfaction degree—depending on its category. The connection is considered as:

- “*unsatisfied*” if the allocated bandwidth is less than its minimum requirement,
- a “*satisfied*” connection if the allocated bandwidth is between its minimum and maximum specified requirements,
- “*over-satisfied*” if it is granted more bandwidth than its maximum need.

The first-tier allocation algorithm is category-based and gives the highest priority to “unsatisfied” connections. For a specific category, the second-tier allocation algorithm is applied to share residual bandwidth based on weights. The flowchart of the proposed 2TSA is shown in Figure 4.7.

Compared to simple-structured scheduling solutions, the hierarchical scheduling mechanisms presented in this section combine in general an inter-service scheduling discipline with a specific queuing mechanism for each service class. Such structures lead to a high computational complexity that may be prohibitive from an implementation point of view and that may not fit the delay constraints of real-time scheduling services.

**Service-specific scheduling** Regardless of the proposed scheduling structure, some service-specific scheduling solutions are presented in literature. Lee *et al.* for instance focused in [43] on VoIP services. They argued that both UGS and rtPS have some problems to support the VoIP services and proposed an enhanced scheduling algorithm to solve the mentioned problems. In fact, the fixed-size grants, assigned to UGS connections of voice users, cause a waste of uplink resources during silence periods. Moreover, the bandwidth request mechanism used by rtPS connections leads to MAC overhead and access delay which is not convenient for VoIP applications. Therefore the authors assumed that a voice activity detector (VAD) or silence detector (SD) is used by the SS in the higher layer and proposed an algorithm to be used by the SSs to inform the BS of their voice state transitions. In order to avoid MAC overhead, the proposed algorithm

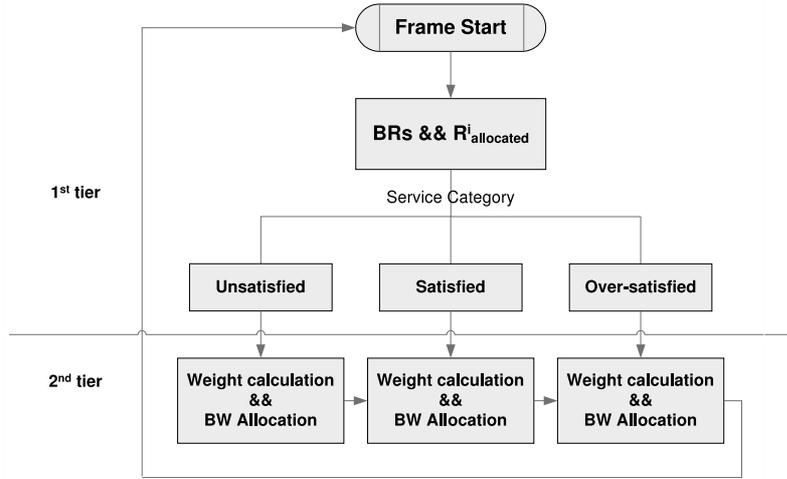


Figure 4.7: Operation flowchart of 2TSA [6]

makes use of one of the reserved bits of the conventional generic MAC header of IEEE 802.16 [1] to do that. Simulation results showed that, compared to rtPS, the proposed algorithm decreases the MAC overhead and access delay. Also it can admit more voice users than UGS making more efficient use of uplink resources.

In a more recent work [44], they demonstrated, using the analysis of resource utilization efficiency, that the ertPS service introduced by the IEEE 802.16e standard [2] is more suitable than UGS and rtPS for VoIP services with variable data rate and silence suppression. Indeed they proved that ertPS not only solves the problems of resource wasting, delay, and overhead caused by the use of UGS and rtPS, respectively but also increases the number of voice users that can be supported by the network.

#### 4.1.2 Optimization-based strategies

This second category of scheduling strategies consists in formulating the scheduling problem, in 802.16 environment, as an optimization problem aiming at optimizing the allocation of resources to different SSs. Table 4.2 presents the formulation of some examples of optimization problems proposed in literature.

To get an optimal solution to the optimization problem formulated in [23] (see Table 4.2), the authors need to use an NP-complete Integer Programming because the number of slots allocated per SS on a given channel should have an integer value. Relaxing this constraint, the authors proposed a second solution based on a linear programming approach that exhibits a complexity of  $O(n^3.m^3.N)$  where  $n$ ,  $m$ , and  $N$  denote the number of SSs, the number of subchannels and the total number of slots, respectively. However, because it is still a computationally demanding problem, the authors suggested the use of a heuristic algorithm whose computational complexity is  $O(n.m.N)$ . The authors then proved that the proposed algorithms optimize the overall system performance but may not be fair to different SSs. Therefore they modified them using the proportional-fair concept.

Based on the developed algorithms, they defined a scheduling algorithm for the BS and another one for the SS. The authors agree that considering a joint scheduling for uplink and downlink, at the BS, is more efficient. They nevertheless argue that it is not possible to do that when con-

sidering the context of OFDMA/TDD. Therefore they adopted a scheduling mechanism in which downlink and uplink are scheduled separately for all the classes. The priorities are assigned as follows. Allocations are made first for UGS, then rtPS, then for nrtPS just to guarantee the minimum requirements, and finally to satisfy the remaining demands. The choice of one of the proposed algorithms depends on the availability of resources and on the channel conditions.

As for the SS, the authors took into account the overall system performance and fairness to different users. They proposed the same sequence followed by the BS but with two different models: a packet model, in which fragmentation is prohibited, for both UGS and rtPS and a byte model—fragmentation is possible—that may be used by nrtPS and BE services.

In [21], Niyato and Hossain considered systems operating in a TDMA/TDD access mode and using WirelessMAN-SC air interface. They defined a utility function that depends on the amount of allocated bandwidth, the average delay, the throughput, and the admission control decision for UGS, rtPS, nrtPS, and BE, respectively. Using these utility functions, they formulated the optimization problem illustrated in Table 4.2. The authors set a limit of the allocated bandwidth between  $b_{min}$  and  $b_{max}$  for each connection. They also defined a threshold for each service class since the total available bandwidth is shared using a threshold-based complete partitioning approach. To obtain the optimal threshold setting, an optimization-based scheme is proposed. To solve the proposed optimization problem, Niyato and Hossain suggested two solutions using an optimal approach and an iterative approach, respectively. The first solution has a complexity of  $O(2^{M(\Delta b)})$  where  $M$  denotes the number of ongoing and incoming connections and  $\Delta b = b_{max} - b_{min} + 1$ . Since the complexity of the optimal algorithm may be prohibitive from an implementation point of view, the authors proposed an iterative approach based the water-filling mechanism. This solution is more implementation-friendly—its complexity is  $O(C)$ —while providing similar system performances.

To analyze the connection-level (such as the blocking probability) and packet-level (e.g. transmission rate) performance measures, the authors developed a queuing and an analytical model, respectively. The proposed connection-level model [21, 22] defines the connection blocking probability and the number of ongoing connections via a Continuous Time Markov Chain (CTMC) model. These parameters are then used to formulate an optimization problem (see Table 4.2) aiming at maximizing the system revenue while maintaining the blocking probability at the target level.

### 4.1.3 Cross-layer strategies

In Sections 4.1.1 and 4.1.2, corresponding to the two first scheduling strategies, we have seen some works (such as [5, 11]) that take into account the AMC capability which is also referred to as MAC-PHY cross layer capability. In those works, the cross-layer aspect is only one of the supported features. However, the scheduling schemes we are presenting in this section are totally found on a cross-layer architecture whose objective is to optimize the communication between different layers of the open systems interconnection (OSI) stack. We can further classify these schemes into: (1) MAC-PHY cross-layer schemes, (2) IP-MAC cross-layer schemes, and application-MAC-PHY cross-layer schemes.

**MAC-PHY cross-layer schemes** The standard provides a link adaptation framework based on which the MCS can be adapted to the channel conditions. However, since no scheduler has been defined by the standard, the way of implementing this capability has been left undefined which

Proposed Solution	Cost Function (Minimize/Maximize)	Constraints (subject to)
Joint Bandwidth Allocation and admission control [22]	Minimize The average delay	<ul style="list-style-type: none"> <li>* The average delay meets the delay requirements of rtPS connections.</li> <li>* The transmission rate meets the transmission rate requirements of connections.</li> <li>* The amount of allocated bandwidth for each connection is between <math>b_{min}</math> and <math>b_{max}</math>.</li> <li>* The total amount of allocated bandwidth does not exceed the total available bandwidth.</li> </ul>
Queuing theoretic and optimization-based model for resource management [21]	Maximize level of users' satisfaction $\Leftrightarrow$ Maximize Utility function	<ul style="list-style-type: none"> <li>* The allocated bandwidth for UGS connections is equal to the required bandwidth</li> <li>* The delay requirements for rtPS connections (depending on the arrival rate, the average SNR and the allocated bandwidth) are met.</li> <li>* The transmission rate requirements of nrtPS connections (depending on the arrival rate, the average SNR and the allocated bandwidth) are met.</li> <li>* BE connections are admitted.</li> <li>* The amount of allocated bandwidth for a given connection is between <math>b_{min}</math> and <math>b_{max}</math>.</li> <li>* The total amount of allocated bandwidth does not exceed the total available bandwidth.</li> <li>* The thresholds (corresponding to the amount of reserved bandwidth for each service class) are respected.</li> </ul>
Queuing model for connection-level performance analysis [21]	Maximize The system revenue $\Leftrightarrow$ Maximize the number of ongoing connections	<ul style="list-style-type: none"> <li>* The connection blocking probabilities* for UGS, rtPS, nrtPS and BE connections do not exceed the target blocking probabilities.</li> </ul>
Efficient and fair Scheduling of Uplink and Downlink in OFDMA Networks [23]	Minimize the unsatisfied demands	<ul style="list-style-type: none"> <li>* The number of granted slots on a given subchannel do not exceed the number of slots of this subchannel</li> <li>* The amount of bandwidth (slots) allocated per connection do not exceed the whole demand of that connection.</li> </ul>

\* The blocking probabilities as well as the number of ongoing connections are function of the corresponding threshold.

Table 4.2: Optimization approach: cost function and constraints

explains the need for such MAC-PHY cross-layer design. This need has been explained and justified through preliminary simulation by Noordin *et al.* in [28] where they propose a cross-layer optimization architecture for WiMAX systems. The cross-layer optimizer (CLO) presented in this work, acts as an interface between between MAC and PHY layers to obtain and tune the required and optimum parameters.

The authors in [28] believe that there is no need to introduce the application layer in the cross-layer architecture they are proposing since the application requirements are considered through QoS provisioning at MAC level. Therefore, the proposed CLO is reduced to MAC-PHY cross-layer optimization.

A more technical MAC-PHY cross-layer scheme has been proposed by Liu *et al.* in [24, 25]. The authors in [24, 25] define an AMC design by setting a region boundary defined by signal to noise ratio (SNR) intervals corresponding each to a different transmission mode. The minimum switching threshold of each interval corresponds to the SNR at which the packet error rate (PER) is less or equal to a prescribed PER  $P_0$ . The AMC design is not adopted for UGS connections because, according to [24, 25], voice traffic can tolerate “some instantaneous packet loss”. Thus, the number of time slots allocated per frame to UGS connections is fixed. Liu *et al.* define a factor called the normalized channel quality based on the received SNR and a priority function (PRF) is assigned to each non-UGS connection depending on its service class. This PRF depends on:

- the BE class coefficient and the normalized channel quality for BE connections,
- the nrtPS class coefficient, the normalized channel quality, and the rate performance for nrtPS connections,
- the rtPS class coefficient, the normalized channel quality, and the delay requirements for rtPS connections.

The class coefficients are set so that the priority order for the different service classes is  $\text{rtPS} > \text{nrtPS} > \text{BE}$ . All the residual time, after scheduling UGS connections, is allocated to the connection having the highest PRF.

The AMC design proposed by Liu *et al.* is quite flexible since it does not depend on any specific traffic or channel model. However, the fact of scheduling only one non-UGS connection per frame might cause a significant delay for real-time applications. This is more likely to happen when the considered PHY is WirelessMAN-OFDM. Indeed, unlike in WirelessMAN-SC PHY where the frame size could take the values: 0.5, 1, or 2 ms, the frame sizes in WirelessMAN-OFDM varies from 2.5 to 20 ms [1, 2] !

Also, the scalability is claimed to be achieved by the proposed scheme since adding new connections would affect connections with low priority prior than those with a high priority. However, this would cause starvation of low priority connections and might even affect high priority ones when the network is overloaded. In order to overcome this shortcoming and guarantee better QoS performance, it would be interesting to combine the proposed scheduling scheme with an efficient CAC algorithm.

**IP-MAC cross-layer schemes** Unlike Noordin *et al.* in [28] who restricted their cross-layer architecture to PHY and MAC layers, the authors in [26, 27] have focused on a layer 3 (L3) and layer 2 (L2) cross-layer design. They insisted on the importance of an IP and MAC cooperation to provide a better QoS service. The cross-layer framework proposed by Mai *et al.* in [26, 27] includes:

- a mapping between L3 and L2 QoS: where integrated service (IntServ) and differentiated service (DiffServ) classes are mapped to 802.16 MAC service classes as shown in Table 4.3.
- a simple admission control scheme based on which a new service flow is accepted when the remaining link capacity is more than the new flow required bandwidth.
- a fragment control mechanism that groups fragments of the same IP packet so that they are treated as a whole by L2 (e.g. fragments from the same IP packet are not interleaved in the L2 buffer, they are all removed in the case of congestion)
- a remapping scheme proposed for a better buffer utilization. Indeed, L3 higher priority CL and EF packets may be stored in nrtPS buffers when rtPS buffers are full (this is more likely to happen because of the burstiness of rtPS traffic).

	IP QoS	MAC 802.16 QoS
<b>IntServ</b>	Guaranteed Service (GS)	UGS
	Controlled load	rtPS
<b>DiffServ</b>	Expedited Forwarding (EF)	nrtPS
	Assured Forwarding (AF)	
<b>IntServ, DiffServ</b>	Best Effort (BE)	BE

Table 4.3: Mapping rule from IP QoS to MAC 802.16 QoS [26, 27]

**Application-MAC-PHY cross-layer schemes** The cross-layer optimization mechanism proposed by Triantafyllopoulou *et al.* in [29, 30] takes advantage of the adaptation capabilities existing at both PHY and application layers. They combine the AMC capability of the physical layer and the multi-rate feature of the multimedia applications through a cross-layer optimizer that exists at BS and SS parts. The optimization process consists in collecting an abstraction of the the layer-specific information (such as QoS parameters and channel conditions) and informing the corresponding layers of the required changes. These changes are instructed based on a decision algorithm that decides about the MCS and traffic rate for each SS.

## 4.2 CAC

In order to guarantee QoS in mobile networks, it is important to combine the scheduling policy with an efficient CAC strategy. The main role of a CAC strategy is to decide whether to accept or not new flows while making sure that the available resources would be sufficient for both the ongoing and the incoming connections. In order to take such an important decision, mainly two strategies can be adopted when no resources are available for the new flows. The first one—more flexible—would consist in gracefully degrading existing connections to make room for the new one. The second strategy—more conservative, yet simpler—would maintain the QoS provided for ongoing connections and simply reject the new service flow.

### 4.2.1 CAC schemes with degradation strategy

This first category of CAC schemes include all the CAC algorithms based on service degradation [45], bandwidth borrowing [46, 47, 48], or bandwidth stealing [49] strategies. The main idea of these policies is to decrease—when necessary and possible—the resources provided to ongoing

connections in order to be able to accept a new service flow. As we will see in this section, this strategy could be combined with a threshold-based capacity sharing approach in order to avoid starvation [49], or a guard channel strategy that reserves a dedicated amount of bandwidth for more bandwidth-sensitive flows (like UGS [47], or handover [48] connections).

#### 4.2.1.1 Service degradation

In [45], service flows (SF) are prioritized according to their respective service type (UGS>(e)rtPS>nrtPS>BE) and among each service type, a priority is assigned to SFs based on their jitter requirements for UGS flows, delay for (e)rtPS flows and traffic priority for both nrtPS and BE flows. If the available bandwidth does not meet the requirements of handover flows, a SF degradation policy is applied. It consists in decreasing the bandwidth assigned to existing SFs whose priority is lower than the handover (HO) SF and whose assigned bandwidth exceeds the minimum reserved bandwidth. SF degradation concerns only handover SFs. A new flow is accepted only if the already available bandwidth guarantees its minimum bandwidth requirement. A two-dimensional continuous Markov model is used to analyze the performance of the proposed scheme. However, many assumptions have been considered: UGS=(e)rtPS and nrtPS=BE. The authors also suppose that all the flow belonging to the same class have the same minimum and maximum requirements which is restrictive. The proposed scheme is then compared to a threshold-based admission control (TAC) policy [15] in terms of blocking and dropping probabilities and bandwidth utilization. Unlike the TAC algorithm, the AC approach proposed by Ge *et al.* [45] adjusts the grant adaptively to the cell load and does not restrict the SF degradation to a single class of flows when necessary. Thus, the proposed algorithm performs better than the TAC algorithm.

#### 4.2.1.2 Bandwidth borrowing

- Bandwidth borrowing in a non-cooperative game

The problem of admission control in IEEE 802.16 networks is formulated by Niyato *et al.* in [46] as a non-cooperative game. The players in this game are the rtPS and nrtPS connections that want to maximize their QoS performance. The payoff of the game is the total utility of the ongoing rtPS and nrtPS connections. The problem consists in finding the equilibrium point between the two types of connections to offer bandwidth for the new connection and meet the QoS requirements of both ongoing and new connection. Based on the solution of the game, a CAC scheme is then proposed to guarantee the QoS requirements of rtPS and nrtPS connections.

- Bandwidth borrowing and stepwise degradation

The CAC scheme, proposed by Wang *et al.* in [47], assigns the highest priority to UGS flows and aims to maximize the bandwidth utilization by bandwidth borrowing and degradation. A predetermined amount of bandwidth  $U$  is exclusively reserved for UGS connections. An UGS connection is accepted if there is enough bandwidth to accommodate its requirements otherwise it is rejected. Denote by  $B$  the total bandwidth, by  $b_{ong}$  the bandwidth set aside for ongoing connections (UGS, rtPS and nrtPS), and by  $b_{ugs}, b_{rtps}$  the bandwidth requirement for a new UGS or rtPS connection, respectively. For a new nrtPS connection,  $b_{nrtps}^{max}$  and  $b_{nrtps}^{min}$  stand for the maximum and minimum bandwidth requirements, respectively. The proposed degradation model is applied when a new rtPS connection is requested and  $b_{ong} + b_{rtps} > B - U$  or when the creation of a new nrtPS connection is requested and  $b_{ong} + b_{nrtps}^{max} - l_{nrtps}^n * \delta \geq B - U$ . where:  $\delta$  is the amount of degraded bandwidth and  $l_{nrtps}^n$  is the current degradation level. Note that only nrtPS

connections could be degraded to accept more rtPS and nrtPS connections. Thus, the reserved bandwidth for each nrtPS connection is  $b_{nrtps}^{max} - l_{nrtps}^n * \delta$  which satisfies  $b_{nrtps}^{max} - l_{nrtps}^n * \delta \geq b_{nrtps}^{max}$  and the maximum degradation level that can be reached is  $(b_{nrtps}^{max} - b_{nrtps}^{min})/\delta$ . In this stepwise degradation scheme, the authors assume that all the connections belonging to the same service type (even non-UGS connections) have the same bandwidth requirements and that the bandwidth requested by an rtPS connection is fixed and does not vary between a maximum sustained and a minimum reserved traffic rates. These assumptions simplify the problem but do not take into account the service requirements specified in the standard.

- Proportional bandwidth borrowing and guard channel

In [48], the authors apply the following priority scheme where handover (HO\_) connections are prioritized over new (N\_) connections: HO\_UGS > HO\_rtps & HO\_ertPS > N\_UGS > N\_rtPS & N\_ertPS > HO\_nrtPS > N\_nrtPS > HO\_BE > N\_BE. The reserved bandwidth corresponds to the maximum sustained traffic rate for UGS and to the minimum required rate for polling services. No bandwidth is reserved for BE traffic. This basic algorithm is combined with a guard channel policy and a proportional bandwidth borrowing scheme. Indeed, a guard channel corresponding to  $n\%$  of the channel capacity is reserved for handover connections. Thus a new connection is blocked if the available bandwidth is less than  $C.n\%$  while a handover connection is blocked only if no bandwidth is available. A proportional bandwidth borrowing scheme is applied when the required bandwidth is not available. The BS borrows from connections having the same or lower priority than the new/HO connection. The connection that occupies more bandwidth lends more to the admitted connection.

#### 4.2.1.3 Bandwidth stealing

In [49], Jiang *et al.* combine an uplink scheduling algorithm with a CAC policy, both based on a token-bucket approach. In the proposed CAC, each uplink connection is characterized by two parameters: a token rate  $r_i$  and a bucket size  $b_i$ . rtPS flows, however, have an extra parameter  $d_i$  corresponding to their delay requirement. In order to avoid starvation of some classes, the authors define a threshold capacity per service type. Thus, a class using more bandwidth than its dedicated threshold has less chances to use the remaining uplink capacity.

When an SS attempts to establish a new service flow—with parameters  $r_i$ ,  $b_i$  and  $d_i$  (for rtPS flows)—with the BS, the proposed CAC algorithm is applied as follows. If the required bandwidth is less than the remaining uplink capacity  $C_{remain}$ , the flow is accepted. If not a "bandwidth stealing" strategy is applied. First, if connections belonging to lower classes—than the new one—are using more bandwidth than their respective thresholds, then the new flow is accepted if the sum of this extra  $C_L$  and  $C_{remain}$  is greater than or equal to its bandwidth requirement. If not, the capacity occupied by connections belonging to the same class of the new one is checked. If it is greater than its threshold, then the new service request is rejected. If not, a bandwidth stealing is attempted from connections belonging to higher classes. This last step is possible only if the capacity of these higher classes exceeds (by  $C_U > 0$ ) their thresholds. If  $C_U + C_L + C_{remain}$  is greater than or equal to the new flow bandwidth requirement, then the new flow is accepted. If not, it is rejected. Note that stealing bandwidth from non-real-time classes (BE and nrtPS) amounts to decreasing their capacity, while for real-time classes it consists in degrading the  $r_i$  of some of their connections to  $c.r_i$  ( $0 < c < 1$ ).

## 4.2.2 CAC schemes without degradation strategy

The hierarchical uplink scheduling algorithm proposed in [13] by Wongthavarawat *et al.* and introduced in Section 4.1.1.2 was combined with a conservative token-bucket-based admission control module. Indeed, no graceful service degradation of existing connections is foreseen by authors to accept a new flow. Thus, a new connection is accepted only if (1) it will receive QoS guarantees in terms of both bandwidth and delay—for real-time flows—and (2) the QoS of existing connections is maintained.

Unlike most of the works where the admission control decision is only based on bandwidth availability, the CAC algorithm proposed by Chandra *et al.* [50] takes also into account the delay and jitter requirements of the service flows. Because the connections have different QoS requirements, an hyper interval (HI) is defined to test the admissibility of the requests. It represents the interval within which the admission process is performed. The authors however consider the delay and jitter requirements for UGS, rtPS and even nrtPS connections which may cause the blocking of an nrtPS connection for instance just because the jitter requirement—which is not necessary in this case as can be seen in Table 3.1—cannot be satisfied. Also, Chandra *et al.* include in their scheme a bandwidth estimator agent that is responsible for monitoring the queue length of both rtPS and nrtPS connections and estimating the bandwidth needs based on the instantaneous change in the queue length. Indeed, the authors define a "configurable threshold"  $BW_{thr}$  according to which, the bandwidth is requested as in the algorithm shown in Figure 6.

---

**Algorithm 6:** Configurable threshold algorithm [50]

---

```

1 Begin
2   if  $((minrate \leq BR) \text{ and } (BR \leq BW_{thr}))$  then
3      $B_{req} \leftarrow minrate$ 
4   else if  $((BW_{thr} \leq BR) \text{ and } (BR \leq maxrate))$  then
5      $B_{req} \leftarrow BR$ 
6   else if  $(maxrate < BR)$  then
7      $B_{req} \leftarrow maxrate$ 

```

---

where:  $BR$  and  $B_{req}$  stand for the bandwidth requirement, and the bandwidth request, respectively.

In [50], the main objective was to ensure QoS guarantee, in terms of bandwidth, delay and jitter. However, only the acceptance ratio was considered to evaluate the performance of the proposed solution.

## 4.2.3 Other CAC schemes

In this section, we introduce some CAC algorithms that have addressed some of the aspects that have not been (or at least not well) investigated in previous works. The first two works [51, 52] have addressed one of the challenges that we have mentioned in Section 3.3 i.e. MAC-PHY cross-layer capability, or more specifically the possibility for a SF to change the burst profile (mainly the MCS)—also known as the AMC capability. We have also chosen to introduce the works done by Yang and Lu in [53, 54] because, unlike the other works presented in previous sections, they have proposed a CAC scheme specifically dedicated for real-time video applications.

---

#### 4.2.3.1 AMC-induced CAC:

[51] is one of the rare works, addressing CAC in 802.16 networks, that take into account the AMC aspect. Indeed, Kwon *et al.* propose an AMC-induced CAC, for IEEE 802.16 networks, that incorporates the modulation type into the CAC process. The work has then been generalized to AMC networks in [55]. The proposed CAC scheme is based on a Markovian model that considers handoff and new connections as well as connections whose modulation changes. The model however supports only two types of modulations and is built based on the assumption that all the connections have fixed and equal bandwidth requirements which limits its applicability.

#### 4.2.3.2 CAC for real-time video applications:

Some CAC solutions existing in literature, have been proposed for a specific kind of applications. In [53] and [54] for instance, the authors have taken advantage of the regularity and periodicity of real-time video traffic to propose a CAC process that particularly fits video applications. Indeed the authors have tried to overcome the time-varying bit rate behavior of video traffics by taking advantage of their group of pictures (GOP) structure—identified by a sequence of I, P and B frames. The main idea consists in avoiding the case where I frames—2 to 10 times bigger than B and P frames—of several flows are transmitted too close to each others. Therefore, the authors have defined a pending period during which the CAC module tries to find a proper time to admit the incoming flow. To fix this proper time, a coordination with I frames algorithm is defined to detect and avoid any I-frame superposition—and thus delay violation—between the ongoing flows and the incoming one. A non-I-frame coordination is then applied. This step aims to place the I and non I frames within their delay bounds. If the CAC is able to perform this step, and this before the pending period expires, the flow is admitted otherwise it is rejected. The amount of data corresponding to non-I frames is computed based on an estimation of non-I-frame rate.

In order to maximize the throughput and minimize the difference of delay between admitted flows, the authors have combined their CAC with a scheduling algorithm. Indeed a latest starting time (LST) algorithm is defined and compared to the EDF algorithm used for instance in [13, 14]. The main limit, which is also the advantage, of this solution is that it only addresses a specific kind of application: real-time video.

Table 4.4 summarizes the different aspects taken into account in the CAC proposals presented in this section. It mainly highlights the criteria (data rate, delay, jitter) based on which the decision, of accepting or rejecting a connection request, has been taken. It also shows whether a degradation and/or a guard channel technique has been adopted by the proposed CAC scheme. Note that we insisted on dedicating a column to AMC even though it has been considered only in [51, 55]. Indeed, we believe that it is a key feature that should not be ignored in the admission control process.

### 4.3 Conclusion

This chapter presents the state of the art of scheduling and CAC algorithms for IEEE 802.16 networks. This survey is by no means an exhaustive compilation of the works addressing this topic. Yet it describes, classifies, and compares scheduling and CAC proposals.

In the last few years, this research area has been intensively investigated and a lot of progress has been done. It is true that CAC and scheduling in wireless networks are classical problems. However, the comparative study presented in this survey shows that, for WiMAX networks, there is still room for improvement.

---

	Data rate	Delay	Jitter	Degradation policy	Guard channel/ Capacity Thresholds	AMC
[45]	✓	—	—	✓	—	—
[47]	—	—	—	✓*	✓**	—
[50]	✓	✓	✓	—	—	—
[48]	✓	—	—	✓	✓***	—
[54, 53] (for video)	✓	✓	—	—	—	—
[49]	✓	✓	—	✓	✓	—
[46]	✓	✓	—	✓	—	—
[51, 55]	✓	—	—	—	✓****	✓

\* stepwise degradation policy, \*\* for UGS connections, \*\*\* for handover connections

\*\*\*\* for handover and modulation changing connections

Table 4.4: CAC in IEEE 802.16 PMP mode: a comparative table

From the scheduling algorithms proposed in literature for IEEE 802.16 networks, we would notice that the main challenging problems that arise when trying to develop a CAC and scheduling strategy are:

- to make a trade-off between an efficient solution, that would take into account the QoS requirements of the different applications, and a simple one that would be implementation-friendly and less time consuming.
- to make a compromise between fairness and channel utilization. Indeed giving priority to users having better channel conditions would increase the channel utilization. Nevertheless, it would be unfair to other users experiencing lower channel conditions.
- to make a choice between an optimized solution that targets a specific kind of applications (like real-time video in [54, 53]) and takes into account its specific needs, and a more general, yet efficient and less complex, scheduling policy that would address heterogeneous types of traffics.
- to take advantage of the adaptive modulation and coding (AMC) capability defined by the standard when proposing a new CAC solution, like it has been proposed in [55]. item to consider the possibility of an adaptive DL/UL bandwidth allocation, as introduced in [7, 9], in order to make an efficient use of the resources and handle unbalanced traffic.

Most of these issues are addressed in our solution described in Chapter 5 where a min-max fairness admission control is adopted and combined to an adaptive DL/UL scheduling algorithm.



## Chapter 5

# Adaptive Scheduling with Max-Min Fairness Admission Control

IEEE 802.16 BWA technology is emerging as a promising solution that provides QoS guarantees for heterogeneous classes of traffic with different QoS requirements. However, despite including the possibility of QoS support, 802.16 MAC protocol does not include a complete solution to offer QoS guarantees for various applications: resource management and scheduling still remain as open issues. In this chapter, we propose a new QoS architecture for PMP 802.16 systems operating in TDD mode over WirelessMAN-OFDM physical layer. It includes a call admission control (CAC) policy and a hierarchical scheduling algorithm. The proposed CAC policy adopts a Min-Max fairness approach making efficient and fair use of the available resources. The proposed scheduling algorithm flexibly adjusts uplink and downlink bandwidth to serve unbalanced traffic. This adaptive per-frame uplink/downlink allocation procedure takes into account the link adaptation capability supported by WiMAX and the data rate constraints of the different types of services. Through simulation, we reveal the efficiency of the proposed CAC scheme and show that our scheduling algorithm can meet the data rate requirements of the scheduling services specified by the IEEE 802.16 Standard. The CAC and scheduling procedures we propose are described in Section 5.1. In Section 5.2, we provide simulation results of our proposal. Finally, Section 5.3 concludes the chapter and gives the possible extensions of the presented work.

### 5.1 Uplink and downlink scheduling

In this section, we present our scheduling proposal. First, we describe the hierarchical scheduling structures proposed for BS and SS. Then, we detail step by step the scheduling algorithm. Finally, we explain the idea of our Min-Max admission control policy.

#### 5.1.1 Hierarchical scheduling structure

As a starting point, we can consider the two-layer hierarchical scheduling structure proposed by Chen *et al* [7]. In first layer scheduling, the authors have suggested two policies. The first one is a transmission direction based priority where they choose to attribute to DL a higher priority than UL. The second policy is a service class based priority applying the following scheme:  $rtPS > nrtPS > BE$ . Additionally, the authors have then combined these policies using a strict priority scheme which assigns strict priority from highest to lowest to:  $DL_{rtPS}, UL_{rtPS}, DL_{nrtPS}, UL_{nrtPS}, DL_{BE}$ ,

---

and  $UL_{BE}$ . For DL and UL UGS connections, they have chosen to apply a fixed bandwidth allocation strategy. In second layer scheduling, they have proposed the use of Deficit Fair Priority Queuing (DFPQ) algorithm. In the scheduling structure we propose, we conserve the hierarchical aspect of scheduling while avoiding the use of cyclic algorithms like DFPQ. We decided to follow this path because in more realistic contexts, the BS does not dispose of enough time to perform such a cyclic scheduling algorithm. We also have a distinct hierarchical organization than the one proposed in [6]. For each level of the hierarchy we decide between:

1. DL and UL: we give a higher priority to downlink for the same reasons given by Chen *et al* [7]. Since we are in the context of a PMP architecture, all the transmissions occur via the BS which is responsible for relaying data between SSs. Also some applications such as HTTP and SMTP require more bandwidth in the downlink.
2. UGS, rtPS, nrtPS and BE: We may combine these two levels as proposed by [7]. Let us nevertheless note that we prefer applying the same scheme for UGS connections. The motivation behind this choice is avoiding resources wasting since it is not necessary to grant a fixed amount of bandwidth for DL UGS connections; the BS is able to adapt the grant to the current needs of each DL UGS connection.
3. Connections having the same scheduling service: since buffers are organized on connection basis, it is important when performing scheduling to know which connection should be served first. Yet, that does not necessarily mean that all the packets of that connection would be served first since other factors such as QoS requirements and availability of bandwidth should be considered too.
  - (a) how to choose between two UGS or two rtPS connections: for both UGS and rtPS scheduling services, we can just adopt a random approach.
  - (b) how to choose between two nrtPS or two BE connections: for these two type of scheduling services, we can take advantage of the Traffic Priority parameter specified in the service flows associated to each nrtPS and BE connection, as shown in Table 3.1.
4. Packets waiting in the same connection queue: not to make the scheduling procedure more complicated, we suggest to use FIFO discipline to schedule packets belonging to the same connection.

As for the SS, the scheduling procedure is easier since the only connections to be managed are those established with the BS in the UL direction. However, as far as the structure is concerned the two scheduling procedures are quite similar. Indeed, the SS follows the following scheme for UL connections:  $UGS > rtPS > nrtPS > BE$ . The scheduling choices we made for the BS at connections and packets levels remain the same for an SS.

Recall that in addition to data transmissions, both BS and SS are asked to schedule MAC management messages. Also the BS is required to poll SSs having at least one (n)rtPS connection, while an SS has to inform the BS of its bandwidth requirements. All these features are considered in our proposal.

### 5.1.2 The BS scheduling algorithm

For the BS, the scope behind the scheduling procedure is to allocate the whole amount of bandwidth available during a frame time interval. Therefore, all the transmissions related to payload,

---

management messages or even gaps and preambles should be elaborately planned and reported in DLFP, UL-MAP and optionally DL-MAP messages. Note that the BS performs the following scheduling procedure at the beginning of each frame interval.

Now, we will describe step by step how this can be possible. We will explain how do the components—at the BS—interact to accomplish the scheduling procedure.

### 5.1.2.1 Step 1: Initialize the available time

Initially, the *BS Scheduler* disposes of a duration equal to the frame time interval—10 ms for example. The scope of this step is first to calculate the time duration corresponding to what is fixed size and should necessarily be sent during the current frame interval, and then to subtract this duration from the frame time interval to know exactly how much time remains. Referring to Figure 2.2, what should be deducted at this level corresponds to the time allocated to TTG and RTG gaps, bandwidth request and initial ranging request intervals, first preamble, and DLFP. Note that all these durations are multiples of the OFDM symbol duration; in other words, no need to perform padding.

So, we update  $T_{av}$  as follows.

$$T_{av} = T_{frame} - \left( T_{pream}^{(long)} + T_{dlfp} + T_{ttg} + T_{opp}^{(rng)} + T_{opp}^{(bw)} + T_{rtg} \right)$$

Note that the notations and parameters used in this chapter are the same reported in Chapter 2.

### 5.1.2.2 Step 2: Plan the first burst

Recall that DL bursts are transmitted in order of decreasing robustness and that the first one contains broadcast MAC control messages. We have specified also that the transmission of DCD and UCD messages is not mandatory unless when at least one of their parameters is updated or when a DCD Interval or an UCD Interval, respectively has elapsed since the transmission of the last DCD or UCD message, respectively.

To make things easier, we suppose that the transmission of two DCD or two UCD messages is spaced by a DCD Interval or UCD Interval, respectively. That is to say that during these intervals the parameters values of DCD and UCD messages are kept unchanged. For DL-MAP message, we can apply the same assumption. Since there is no need to send a DL-MAP message unless there are more than four DL bursts—corresponding to four different burst profiles—or unless a Lost DL-MAP Interval has elapsed since the transmission of the last DL-MAP message. For sake of simplicity, we assume that we do not exceed four burst profiles on DL. In this case, only DLFP is needed to describe the location and profile of DL bursts and one full DL-MAP must be broadcast in the first burst within the Lost DL-MAP Interval, as it is specified in [1].

Once the first burst is planned,  $T_{av}$  parameter should be updated as follows (c.f. Chapter 2):

$$T_{av} = T_{av} - \frac{L_{bst}[1] + L_{pad}[1]}{L_{sym}[1]} * T_{sym}$$

Note that  $T_{av}$  must always be a multiple of an OFDM symbol duration.

### 5.1.2.3 Step 3: Proceed in accordance with the scheduling structure

After calculating the time it disposes of for DL and UL transmissions, the *BS Scheduler* performs scheduling taking into account the hierarchical scheduling structure described in Section 5.1.1. Indeed the first service type to consider is UGS and more specifically the UGS DL connections. For each scheduling service the *DL Grant Allocator* or *UL Grant Allocator*—depending on the connection direction—proceed as follows:

**Determine the number of packets to serve per connection** In order to determine this number, we need the following parameters:

- $S_{gmh}$ : size (in bytes) of a generic MAC header.
- $S_{brh}$ : size (in bytes) of a bandwidth request header.
- $S_{crc}$ : size (in bytes) of a CRC field.
- $S_{ulmap\_ie}$ : size (in bytes) of an UL-MAP IE.

Consider the following parameters associated to a given connection  $j$ :

- $n_i^j$ : the number of packets of connection  $j$  that are transmitted during the  $i^{th}$  frame interval.
- $N_i^j$ : the number of packets of connection  $j$  that are transmitted during the  $i$  last frame intervals.
- $R_{max}^j$ : the Maximum Sustained Traffic Rate of connection  $j$ .
- $R_{min}^j$ : the Minimum Reserved Traffic Rate of connection  $j$ .
- $R^j$ : the rate to be considered during the scheduling procedure; with  $R_{min}^j \leq R^j \leq R_{max}^j$ .  $R^j$  corresponds to the maximum actual rate at which the connection may be allowed to transmit its data. This parameter is calculated by the *Admission Control Module* in such a manner that allowing the considered connection to transmit at this rate would not affect the QoS of existing connections; For UGS and BE connections  $R^j = R_{max}^j$ ; further details on how this parameter is computed are given in Section 5.1.3.
- $R_i^j$ : the amount of requested bandwidth (bits) for connection  $j$ . This request is sent during the  $(i-1)^{th}$  frame interval in order to be satisfied during the  $i^{th}$  frame interval. This parameter includes payload and MAC overhead but not physical one. The use of this parameter is meaningless in the case of a UGS connection.
- $Q_i^j$ : the number of packets that are waiting in the queue of connection  $j$ . This parameter concerns only DL connections.
- $S_{pkt}$ : packet size (in bytes).

In the beginning of each frame interval  $i$ , the number of packets to transmit per connection should be calculated given  $N_{i-1}^j$ ,  $R^j$ ,  $S_{pkt}$ , and possibly  $Q_i^j$  when  $j$  is a DL connection or  $R_i^j$  when  $j$  is a non-UGS UL connection.

To compute  $n_i^j$ , we shall consider the three following cases:

- **Case 1:  $j$  is a DL connection**

The idea is that the *DL Grant Allocator* tries to offer to connection  $j$  the possibility of transmitting a number  $n_i^j$  of packets big enough to guarantee for connection  $j$  reaching the maximum rate allowed by the *Admission Control Module*. Of course  $n_i^j$  could not exceed the number of packets waiting in the queue of connection  $j$ . Note that this is applied even for UGS DL connections in order to avoid potential bandwidth wasting.

$$n_i^j = \min \left( \left\lceil \frac{R^j * i * T_{frame}}{S_{pkt} * 8} \right\rceil - N_{i-1}^j, Q_i^j \right) \quad (5.1)$$

- **Case 2:  $j$  is an UL UGS connection**

Since we are considering an UGS connection—in which case the Maximum Sustained Traffic Rate corresponds to the Minimum Reserved Traffic Rate and also to the maximum rate allowed by the *Admission Control Module*—the *UL Grant Allocator* should offer to connection  $j$  the possibility of transmitting a number  $n_i^j$  of packets big enough to guarantee to connection  $j$  to reach the Maximum Sustained Traffic Rate specified in its service flow:

$$n_i^j = \left\lceil \frac{R_{max}^j * i * T_{frame}}{S_{pkt} * 8} \right\rceil - N_{i-1}^j \quad (5.2)$$

- **Case 3:  $j$  is a non-UGS UL connection**

Unlike the previous case, the bandwidth requirements of non-UGS connections must be explicitly formulated by the SS Scheduler—more specifically by the *BW REQ Manager*—which has a more accurate perception of the UL queues status. This bandwidth request, corresponding here to the parameter  $Req_i^j$ , is formulated during the  $(i-1)^{th}$  frame interval and represents the amount of bandwidth needed during the  $i^{th}$  frame interval. However, the *BS Scheduler* must check whether  $Req_i^j$  exceeds what has been fixed by the *Admission Control Module*; in which case, the *UL Grant Allocator* performs shaping by choosing the minimum between what has been requested by the *SS Scheduler* and what would normally be planned by the *BS Scheduler* in order to guarantee the maximum rate allowed by the *Admission Control Module* for connection  $j$  for the  $i$  last frame intervals:

$$n_i^j = \min \left( \frac{Req_i^j}{S_{gmh} + S_{pkt} + S_{crc}}, \left\lceil \frac{R^j * i * T_{frame}}{S_{pkt} * 8} \right\rceil - N_{i-1}^j \right) \quad (5.3)$$

**Calculate the resulting overhead and check the availability of bandwidth** The scope of this step is to calculate the overhead that would result from the transmission of  $n_i^j$  packets of a given connection  $j$  and then to check if the remaining bandwidth allows such transmission. Suppose that the transmission must occur during burst  $k$ .

Let us consider the following variables:

- $tmp\_L_{bst}[k]$ : a temporary variable used to estimate the value of  $L_{bst}[k]$  if the  $n_i^j$  packets are transmitted.
- $tmp\_L_{pad}[k]$ : a temporary variable used to estimate the value of  $L_{pad}[k]$  if the  $n_i^j$  packets are transmitted.

- $tmp\_T_{bst}[k]$ : a temporary variable used to estimate the duration of burst  $k$  if the  $n_i^j$  packets are transmitted.

To calculate the overhead resulting from the transmission of  $n_i^j$  packets, we should consider the following two cases:

- **Case 1:  $j$  is an UL connection whose SS has not received any grant during the current frame**

In this first case, an UL-MAP IE should be addressed to the considered SS since receives an UL grant for the first time during the current frame. Yet, adding an IE to the UL-MAP message would impact the length and the number of padding bits of the first burst as follows:

$$tmp\_L_{bst}[1] = L_{bst}[1] + S_{ulmap\_ie}$$

$$tmp\_L_{pad}[1] = compute\_pad(tmp\_L_{bst}[1], 1)$$

$compute\_pad()$  is a function that returns the size of the necessary padding when given the payload size of a burst and its index—to get the number of bits per symbol associated to that burst.

As it is the first UL grant to be addressed to the considered SS, we should add a short preamble to burst  $k$  for PHY synchronization (see Figure 1.1). Yet, since a short preamble duration corresponds to the duration of one OFDM symbol, we can just add the number of bits per symbol to the corresponding burst length:

$$tmp\_L_{bst}[k] = L_{bst}[k] + L_{sym}[k] + n_i^j * (S_{gmh} + S_{pkt} + S_{crc}) * 8$$

The duration of burst  $k$  is given by:

$$tmp\_T_{bst}[k] = \frac{tmp\_L_{bst}[k] + tmp\_L_{pad}[k]}{L_{sym}[k]} * T_{sym}$$

Once the overhead that may be introduced by the transmission of  $n_i^j$  packets is calculated and once the duration of the associated burst is known, the *UL Grant Allocator* executes Algorithm 7 to check whether it is possible to send these  $n_i^j$  packets while taking into account the remaining time.

---

**Algorithm 7:** Compute overhead: Case 1

---

```

1 Begin
2   if  $((T_{av} + T_{bst}[k] - tmp\_T_{bst}[k] + T_{bst}[1] - tmp\_T_{bst}[1]) > 0)$  then
3      $L_{pad}[k] \leftarrow tmp\_L_{pad}[k]$ 
4      $L_{bst}[1] \leftarrow tmp\_L_{bst}[1]$ 
5      $L_{pad}[1] \leftarrow tmp\_L_{pad}[1]$ 
6      $N_i \leftarrow N_{i-1}^j + n_i^j$ 
7      $T_{av} \leftarrow T_{av} + T_{bst}[k] - tmp\_T_{bst}[k] + T_{bst}[1] - tmp\_T_{bst}[1]$ 

```

---

- **Case 2: j is a DL connection<sup>1</sup> or an UL one whose SS has received a grant in the current frame**

In comparison to what was specified in the first case, there is no need in this case to add neither a preamble nor an IE in the UL-MAP message.

$$tmp\_L_{bst}[k] = L_{bst}[k] + n_i^j * (S_{gmh} + S_{pkt} + S_{crc}) * 8$$

$$tmp\_L_{pad}[k] = compute\_pad(tmp\_L_{bst}[k], k)$$

Like in the first considered case, after estimating the overhead that would result from the transmission of  $n_i^j$  packets, the *UL/DL Grant Allocator*—depending on the connection direction—executes Algorithm 8 to check the remaining time makes it is possible to plan the transmission of these  $n_i^j$  packets.

---

**Algorithm 8:** Compute overhead: Case 2

---

```

1 Begin
2   if  $((T_{av} + T_{bst}[k] - tmp\_T_{bst}[k]) > 0)$  then
3      $L_{bst}[k] \leftarrow tmp\_L_{bst}[k]$ 
4      $L_{pad}[k] \leftarrow tmp\_L_{pad}[k]$ 
5      $N_i^j \leftarrow N_{i-1}^j + n_i^j$ 
6      $T_{av} \leftarrow T_{av} + T_{bst}[k] - tmp\_T_{bst}[k]$ 

```

---

#### 5.1.2.4 Step 4: Share bandwidth and plan transmissions

Once the number of packets to be scheduled and the resulting overhead are determined, the *BS Scheduler* plans the data transmissions as well as DSx messages and polling opportunities decided by the *DSx Manager* and the *Polling Manager*, respectively. It then equally shares the remaining bandwidth (if any) among the SSs in terms of OFDM symbols in order to avoid padding. SSs having only UGS connections are not concerned by this grant since their needs are wholly satisfied. Based on all these scheduling decisions, the *BS Scheduler* generates the DLFP and UL-MAP messages and broadcast them on the downlink.

### 5.1.3 Admission control policy

The purpose behind adopting an admission control policy is to satisfy the QoS requirements of new service flows while respecting the QoS constraints of existing connections and trying to be as fair as possible when granting resources. In order to simplify the admission control algorithm, we will only consider active connections; so no bandwidth will be reserved or granted to preprovisioned or admitted service flows.

As mentioned in Section 5.1.2, the rate to be considered in the scheduling procedure—and which corresponds to the maximum rate allowed for a given connection—is computed during the admission control time period. This rate is determined in such a manner to:

---

<sup>1</sup>Since the DLFP message has a fixed size, addressing a grant to a DL connection does not imply any modification for the first DL burst.

---

- guarantee at least the Minimum Reserved Traffic Rate for all accepted connections;
- make efficient use of current available resources; this use should be adapted to channel conditions;
- try to be as fair as possible when it is necessary to degrade or ameliorate the QoS of existing active flows.

Based on these criteria, we define first when the admission control procedure should be applied and then how the maximum rate to allow for each rtPS and nrtPS connections may be computed. Actually, the admission control mechanism is performed when an SS or a BS attempts to establish a new active connection and also when an SS uses either a more or a less robust DL or UL burst profile. The motivation behind considering these two latter cases may be explained as follows. When an SS uses a more efficient burst profile, this means that it will need less time and thus less resources to keep the same rates. The admission control policy should then make use of this extra bandwidth and try to share it among existing active flows and hence improve their respective maximum allowed rate, corresponding to parameter  $R^j$  (see Section 5.1.2). In this case the purpose behind applying the admission control mechanism is not accepting or rejecting a request but updating the  $R^j$  parameter (of polling connections) to be used during the scheduling procedure. This corresponds also to the case when an active connection is deleted. When an SS should use a more robust MCS, this may affect existing connections since the SS would need more resources to keep the same rates. Therefore, the admission control mechanism should redistribute the available resources (recompute  $R^j$ ) and reject if necessary one or more connections.

The admission control mechanism proceeds as follows:

- It accepts all BE addition requests since they don't have any QoS requirements.
- It checks whether it is possible to guarantee the Maximum Sustained Traffic Rate for all the considered non-BE connections: existing connections and the one attempting to be established. To do that, :
  - it first computes the ceiling number  $n$  of packets to serve per frame for each connection. This number is computed based on the Maximum Reserved Traffic Rate specified in the SF associated to connection  $j$ . The number of packets is given by:

$$n = \left\lceil \frac{R_{max}^j * i * T_{frame}}{S_{pkt} * 8} \right\rceil$$

- then it calculates the overhead resulting from the transmission of these  $n$  packets based on the same approach applied for scheduling. Nevertheless, in order to facilitate this step, we assume that all the active SSs (having at least one connection) belonging to the network receive grant in each frame. In other words, the UL-MAP message to consider when computing the resulting overhead contains as much Data Grant IEs as the number of active SSs. This implies a fixed-size UL-MAP message and then a fixed-size part for MAC management messages; we also consider the worst case corresponding to the case where full DCD, UCD and DL-MAP messages are sent in the considered frame. Furthermore, if  $j$  corresponds to an rtPS connection or an nrtPS, we should also take into account the amount of bandwidth necessary to poll the associated SS. The polling period to consider here is the same considered by the *Polling Manager*. If the available bandwidth allows such grants—for all the considered connections—then the new connection is accepted and  $R^j$  corresponds to  $R_{max}^j$ , otherwise:

	Connection parameters			Simulation Time (frame interval)		
	$R_{max}^j$ (Mbps)	$R_{min}^j$ (Mbps)	$S_{pkt}$ (B)	0-500 UIUC = QPSK1/2	500-2000	2000-5000 UIUC = 16-QAM 3/4
DL UGS	6	6	1500	✓	✓	✓
UL UGS	2	2	1500		✓	✓
UL rtPS #1	4	2	1500	✓	✓	✓
UL rtPS #2	5	1	1500	✓	✓	✓
UL BE	1	0	1500	✓	✓	✓

Table 5.1: Single SS Scenario parameters

- \* It checks whether it is possible to guarantee the Minimum Reserved Traffic Rate for all the considered non-BE connections. This step is performed similarly to the previous one, just replacing  $R_{max}^j$  by  $R_{min}^j$ . If it is not possible, the connection addition request is rejected, otherwise  $R^j$  is set to  $R_{min}^j$  and:
- \* If there is a remaining amount of bandwidth, it is shared among existing rtPS and nrtPS connections since only these services have specific QoS requirements and may have better rates than the Minimum Reserved Traffic Rate. In order to avoid padding, the sharing will be made in terms of OFDM symbols. Moreover, trying to be as fair as possible, we adopt a Min-Max weighted fair allocation to share the remaining bandwidth. The proposed Min-Max mechanism considers the channel conditions experienced by each SS—associated SFs—by assigning weights inversely proportional to the efficiency (number of bits per symbol) of the corresponding MCS. Each (n)rtPS SF is then allocated a percentage of available OFDM symbols based on a normalized weight. Further, the remaining bandwidth is redistributed among unsatisfied (n)rtPS flows according to their new normalized weights. The process continues till no bandwidth is left.

The graph in Figure 5.1 illustrates the admission control policy explained above.

## 5.2 Performance analysis

In this section, we present two simulation scenarios to study the efficiency and fairness of the proposed admission control and scheduling solution. The goodput is the main parameter targeted in these scenarios; it is studied under different conditions and network configurations. Simulations were carried on MATLAB and address systems operating in TDD mode over WirelessMAN-OFDM physical layer. In the considered scenarios, we study the performance of the proposed CAC and scheduling solution for different scheduling services in a network involving one or multiple SSs using the same or different MCSs. In order to test to which extent would the system adapt its grants to the network conditions, we have split the scenarios into different configuration phases each corresponding to a specific set of connections and channel conditions.

### 5.2.1 Scenario 1: Single SS scenario

In this scenario, we consider a basic network structure composed of one BS and only one SS. During the first interval, corresponding to 500 time frame intervals, the SS establishes with the BS one UGS connection, two rtPS connections and one BE connection. Table 5.1 presents the main parameters of each connection: the maximum sustained traffic rate, the minimum reserved traffic

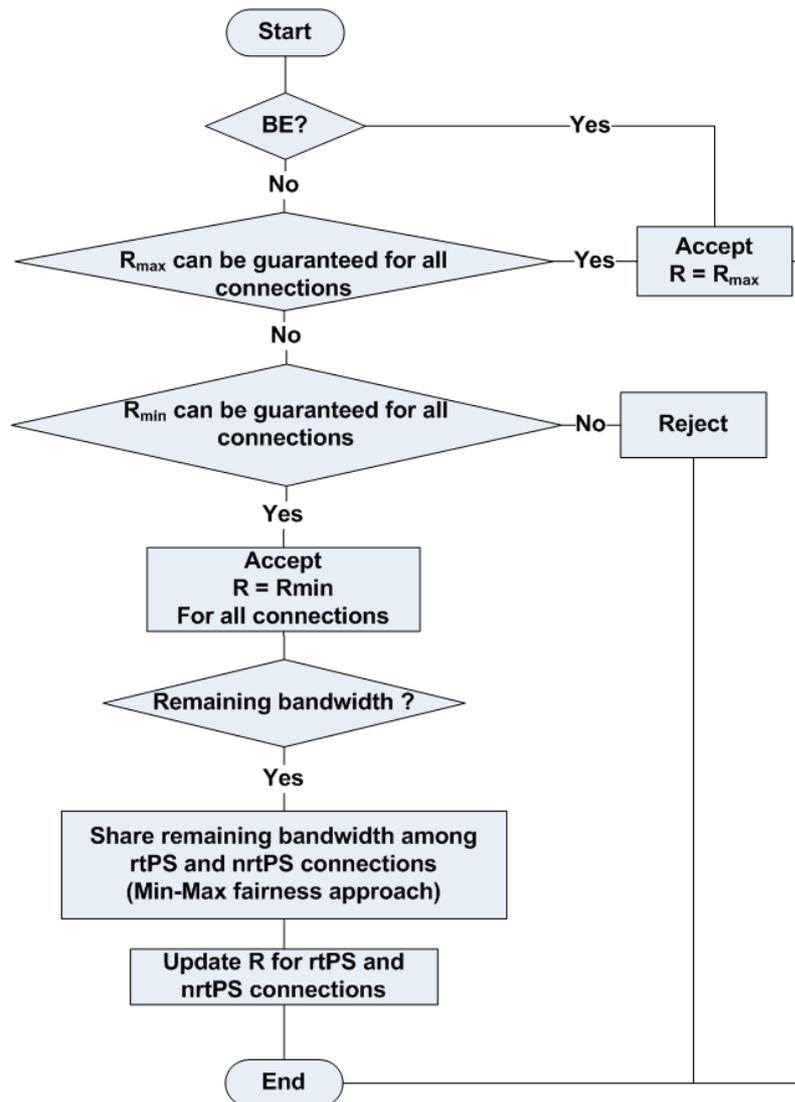


Figure 5.1: Min-Max CAC Policy

rate and the packet size. The MCS used by the SS during this first time interval is QPSK 1/2. During the second time interval (see Table 5.1), the SS establishes a new UGS connection with the BS. The reason to introduce a new UGS connection is to test the efficiency and the fairness of the proposed CAC module and scheduling algorithm and to see what would be the impact of such event on the QoS of existing connections. During the last time interval considered in this scenario (from the 2000<sup>th</sup> to the 5000<sup>th</sup> frame), we suppose that the SS has better channel conditions and may use a more efficient MCS (16-QAM 3/4). The motivation behind this is to test how would the BS and SS Schedulers adapt their grants to the new channel conditions and take advantage of this extra bandwidth. Note that we have chosen to generate the traffic at the maximum rate for the three services in order to make sure that obtained results are due to the scheduling policy and not to the way the traffic is generated.

Figure 5.2 presents the goodput for each connection. Let us focus first on the 1<sup>st</sup> to 500<sup>th</sup> frame interval. As shown in Figure 5.2, the four connections configured in this interval are accepted by the CAC module; the UGS connection is granted the maximum (which corresponds also to the

minimum) traffic rate specified in its service flow. The two rtPS connections share the remaining of the bandwidth quite fairly considering the minimum rate specified by each of them. However, the BE connection suffers from starvation since it has no minimum QoS requirements compared to the other services. As we can see in Figure 5.2, the new UGS connection established on the 500<sup>th</sup> frame interval is accepted by the CAC module and is offered the required rate specified in its service flow. This causes the degradation of both rtPS connections with approximately the same rate while providing them nevertheless with more than the minimum rate specified in their respective service flows. However, this allows the BE connection to obtain some grant. The fact that the BE connection succeeds in having some resources in this interval may be explained as follows. During the first interval, the amount of bandwidth remaining after providing the two rtPS connections with the rate calculated via the Max-Min algorithm was not enough to send a packet from the BE backlog. During the second interval however, when the rates of rtPS connections have decreased, the reached rates values let sufficient extra bandwidth to send at least one packet from the BE connection queue. So the packet-based scheduling policy we use may decrease in some cases the risk of starvation for BE connections. Let us consider now the last interval during which the channel experiences better conditions. Figure 5.2 shows that the scheduling procedure we propose dynamically adapts the grants to the new channel conditions and makes efficient use of the new available resources. Indeed, rtPS connections as well as BE connection succeed to get the maximum rate specified in their SFs.

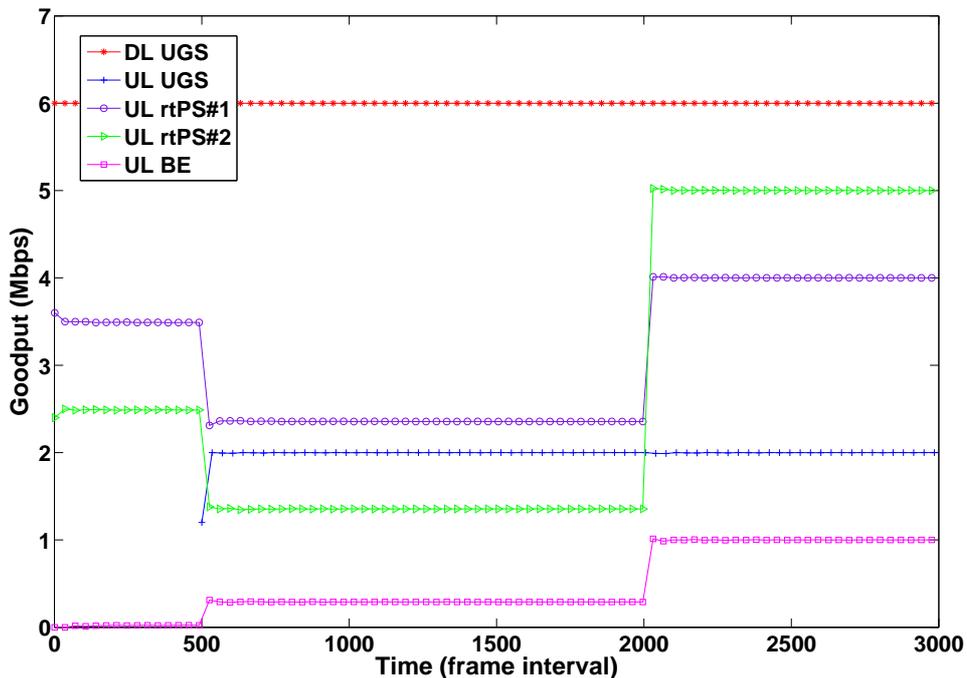


Figure 5.2: Single SS scenario

### 5.2.2 Scenario 4: Multiple SSs scenario

The network considered in this scenario is composed of one BS and three SSs. We assume that the three SSs have the same channel conditions and use the same MCS. They enter the network at

	SF parameters			Simulation Time (frame interval)		
	$R_{max}^j$ (Mbps)	$R_{min}^j$ (Mbps)	$S_{pkt}$ (B)	0-1000	1000-4000	4000-8000
UGS	5	5	1500	✓	✓	✓
rtPS	6	2	1500	✓	✓	✓
nrtPS	3	1	1500	✓	✓	✓
BE	0.7	0	1500	✓	✓	✓
UGS	2.5	2.5	1500		✓	✓
rtPS	5	3	1500		✓	✓
nrtPS	7	2	1500		✓	✓
BE	0.3	0	1500		✓	✓
UGS#1	1	1	1500			✓
UGS#2	3	3	1500			✓
nrtPS#1	1.5	0.5	1500			✓
nrtPS#2	2.5	0.5	1500			✓

Table 5.2: Multiple SSs parameters

different time intervals. Each of them tries to establish a set of connections with the BS and vice versa. The parameters of these connections as well as their durations are reported in Table 5.2.

Figure 5.3 depicts the goodput of each connection during the simulation time. In the first time interval (from frame 1 to frame 1000), only SS1 is connected to the BS. Its connections are granted the maximum rate specified in their respective service flows. When SS2 enters the network, we note that SS1 UGS connection conserves the same rate while SS1 rtPS and nrtPS connections goodputs decrease and reach values a little bit more than their respective minimum reserved rates. Both BE connections succeed however to have the maximum sustained traffic rate specified in their SFs. At frame 4000 (corresponding to 40s after the beginning of the simulation), SS3 joins the network. Only two of the four connections it attempts to establish with the BS are accepted by the CAC module. The network has reached its maximum capacity. Indeed, all connections are getting granted only the minimum reserved rate and this explains the CAC decision since accepting another connection would have degraded the QoS of existing ones. After granting UGS and polling connections, the remaining bandwidth allows nevertheless SS1 and SS2 BE connections to send some data.

### 5.3 Conclusion

In this chapter, we have proposed a new adaptive QoS architecture for PMP 802.16 systems operating in TDD mode over WirelessMAN-OFDM physical layer. The proposed architecture includes a CAC module and a hierarchical scheduling structure. The CAC module we have proposed flexibly adjusts the grants boundaries to the connections QoS requirements while making efficient and fair use of the dynamic channel capacity via a Min-Max fairness approach. The proposed scheduling procedure adapts the frame-by-frame allocations to the current needs of the connections with respect to the grants boundaries fixed by the CAC module. These boundaries may be set through a degradation of the ongoing connections rates if the available resources are not enough to accommodate the needs of a new connection for example. This degradation can be handled by UDP traffic. However, it might cause an uneven behavior for TCP traffic especially under short round trip time (RTT) conditions. To prevent such a behavior, we can combine our CAC policy

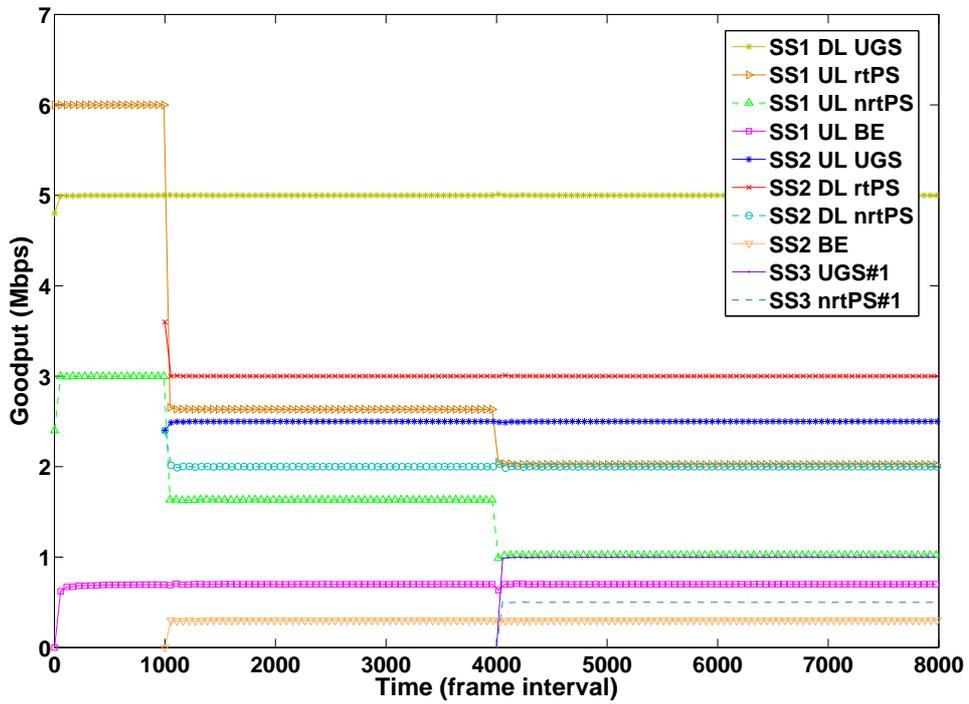


Figure 5.3: Multiple SSs Scenario

with a TCP-friendly traffic policing mechanism among those available in the literature [56]. A further challenge we face will be to support bursty traffics and to integrate delay constraints in our proposal. These two shortcomings are addressed in our multi-Constraints Scheduling Strategy (mCoSS) which is presented in next chapter.



## Chapter 6

# mCoSS: a multi-Constraints Scheduling Strategy for WiMAX Networks

In this chapter, we attempt to assemble the different pieces of the resource allocation puzzle of mobile WiMAX networks by addressing the main scheduling issues that are still open. We thus propose a multi-Constraints Scheduling Strategy (mCoSS) which specifies the scheduling-related operations at both the base stations and the mobile stations. The proposed scheduling strategy is described through a set of scheduling algorithms that maximize the quality of service (QoS) degree of satisfaction for both real-time and non-real-time traffic in terms of delay and throughput. The proposed strategy can be applied to both OFDM and band-AMC OFDMA environments. The access to the network is regulated via a traffic shaper that is inspired from the dual token-buckets shaping mechanism which allows traffic burstiness while protecting contract-conforming connections from misbehaving ones. The dual-bucket mechanism is combined with a two-rounds scheduling algorithm reflecting the two levels of service to be expected by each connection. In the first round, the minimum reserved traffic rate and delay constraints are met while in the second round, fairness among flows is ensured over the remaining bandwidth due to a weighted fair queuing (WFQ) mechanism. The bandwidth request and grant policy adopted in the proposed strategy takes advantage of the different mechanisms specified by the IEEE 802.16e standard and adapts the choice of the appropriate technique to the service flow QoS constraints and to the current availability of radio resources. Other concerns such as supporting the link adaptation capability and avoiding starvation of best effort traffic are also considered. To evaluate the performance of mCoSS, we have implemented the corresponding set of algorithms in QualNet simulator and compared them to strict priority (SP) and to a variant of WFQ discipline. The obtained results show a nice tradeoff between fairness and efficiency with a high respect for the connections' QoS requirements.

The remainder of this chapter is organized as follows. Section 6.1 explains the idea of the modified dual-bucket traffic shaping mechanism adopted in our strategy. In Section 6.2, we provide the details of the proposed two-rounds scheduling approach used by the MS and BS for DL and UL. The proposed bandwidth request and grant policy is described in the same section. The performance evaluation of mCoSS is given in Section 6.3 after describing the OFDMA-based WiMAX simulation model provided by QualNet. Section 6.4 concludes the chapter by summarizing the main features supported by mCoSS and pointing out the main obtained results.

---

## 6.1 A modified dual-bucket shaping mechanism

In order to provide QoS for different types of flows, it is important to implement a traffic shaping mechanism to control the volume of traffic entering the network and to isolate well-behaving traffics from misbehaving ones. The two main traffic shapers implementations used in traffic engineering are: the leaky bucket and the token bucket. The leaky bucket provides a mechanism by which a flow is shaped to be sent to the network at a constant rate. The token bucket however, while providing rate control, allows the traffic to burst up to a configurable threshold. In order to accommodate the bursty characteristics of some categories of applications targeted by WiMAX, we choose the latter mechanism to model our traffic shaper. More specifically, we use the multiple-buckets variant of the token-bucket implementation. We associate each flow  $i$  with two buckets corresponding to the minimum reserved traffic rate  $R_{min}^i$ , and to the maximum sustained traffic rate  $R_{max}^i$ . These per-flow dual buckets reflect the lower and upper boundaries of the service to be provided for each flow. Each bucket has three components: a burst size, a mean rate and a time interval. Figure 6.1 represents the dual-buckets structure associated to a service flow. The first bucket is characterized by:

- a mean rate, also called committed information rate ( $CIR$ ), which specifies the amount of data that can be sent per time unit on average.
- a time interval  $T_c$ , also called the measurement interval; it specifies the time quantum in second per burst.
- a burst size, also called committed burst size ( $B_c$ ); it corresponds to how much traffic can be sent per burst within a given measurement interval.

The three parameters are linked as follows:  $CIR = \frac{B_c}{T_c}$ . We set  $CIR$  to the minimum reserved traffic rate  $R_{min}^i$ , and  $T_c$  to a grant interval  $I_{gr}^i$  characterizing the  $i$  flow. For a real-time traffic  $i$ , this parameter corresponds to the maximum latency  $L_{max}^i$ . For non-real time flows, this parameter should not exceed the polling interval (for nrtPS) and might be set to a value that is a function of the mean transmitting interval of the flow. The introduction of this parameter is needed first to define the frequency of the allocations for each flow and because the standard does not specify the interval over which  $R_{min}^i$  and  $R_{max}^i$  are averaged. This first bucket reflects basically the service level agreement (SLA) a WiMAX system is committed to provide for a flow. Recall that, as mentioned in Section 3.1, a BS or SS does not have to meet the latency service commitment ( $L_{max}$ ) for service flows that exceed their minimum reserved rate [37].

The second bucket is used to make sure that the rate at which the traffic is transmitted stays within the allowed boundaries; i.e. it does not exceed  $R_{max}^i$ . As shown in Figure 6.1, the second bucket is defined through the following components: a mean rate called excess information rate ( $EIR$ ), an excess burst size  $B_e$ , and a time interval  $T_e$ . In order to average the rate over the grant interval of the flow, we consider the same measurement interval. i.e.  $T_e = T_c = I_{gr}^i$ . More specifically, for a real-time flow  $i$ ,  $T_e = T_c = L_{max}^i$ .  $B_e$  is configured in such a way that the maximum burst size does not exceed  $R_{max}^i \times T_e$ . In other words,  $B_c + B_e = R_{max}^i \times T_e$  which implies that  $B_e = EIR \times T_e = (R_{max}^i - R_{min}^i) \times T_e$ . Note that when the capacity of the buckets  $B_c$  or  $B_e$  is reached, all the extra tokens are discarded. Using the configuration described above, if the buckets are empty at the beginning of the grant interval, the maximum burst size can be only reached at the end of the grant interval if no tokens are removed meanwhile. More specifically, if the packets are generated at  $R_{max}^i$  in a bursty way (still contract-conforming), they need to be delayed even if there are enough resources to transmit them since there are no enough tokens in the buckets. This configuration allows to smooth the traffic and to avoid bottlenecks at the next hop.

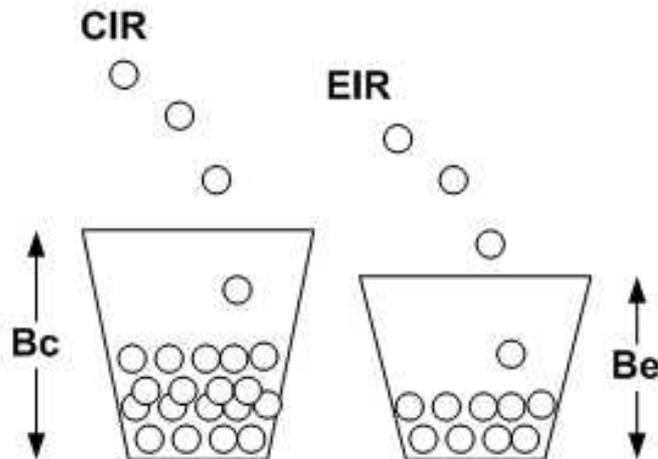


Figure 6.1: A Dual-Bucket Shaping Mechanism

Nevertheless, it might lead to a waste of resources. For more flexibility in resource management and in order to reach a better frame utilization rate, we choose to implement a modified version of the dual-bucket mechanism previously described. In this modified configuration, we keep the same values of the measurement intervals  $T_c$  and  $T_e$ , and burst sizes  $B_c$  and  $B_e$ . Nevertheless, we consider the buckets full at the beginning of the interval. This configuration, while bounding the burstiness to the allowed thresholds, allows it to occur at anytime during the grant interval. Note that for BE connections, the first bucket is empty since  $CIR = R_{min}^i = 0$  and for UGS connections the second bucket is empty since  $R_{max}^i = R_{min}^i$  and  $EIR = R_{max}^i - R_{min}^i$ . Thus, the same settings remain applicable to all scheduling service types. The proposed traffic shaping is combined with a two-rounds scheduling algorithm. More details about the whole mechanism are provided in next section.

## 6.2 A two-rounds scheduling algorithm

The scheduling framework we propose in this chapter consists of three schedulers; two running at the BS: one for DL and one for UL and a scheduler running at the SS to redistribute the bandwidth allocated by the BS among the UL connections. Moreover, the UL schedulers (both at the BS and the SS) rely on a bandwidth request and grant process that allows the SS to transmit its non-UGS bandwidth needs to the BS which would decide the bandwidth grants accordingly. In this Section, the three scheduling processes are described. At the beginning of each frame, the BS has to decide about the way of sharing the available bandwidth among active service flows. The scheduling process we propose consists of two scheduling rounds.

During the first round of the scheduling process, the objective is to honor the SLA by providing the minimum reserved traffic rate to non-BE active connections and by meeting the latency requirements of real-time services (UGS, ertPS, and rtPS). The frequency of these first allocations is set to the scheduling grant interval of the flow:  $I_{gr}^i$ . Referring to the dual-bucket mechanism described in the previous section, this first scheduling round is aimed at emptying the first token bucket of the flows whose grant interval expires in current frame interval. By proceeding this way, we avoid to schedule every single connection at each frame interval which decreases the overhead associated to a per-SS access. The algorithms corresponding to the implementation of this first round at the BS (DL and UL) and at the SS are provided in Algorithm 9, Algorithm 11, and

Algorithm 10, respectively. The parameters considered in these algorithms are the following:

- $U = \{u1, u2, \dots, uu\}$  the set of UGS SFs
- $E = \{e1, e2, \dots, ee\}$  the set of ertPS SFs
- $R = \{r1, r2, \dots, rr\}$  the set of rtPS SFs
- $N = \{n1, n2, \dots, nn\}$  the set of nrtPS SFs
- $B = \{b1, b2, \dots, bb\}$  the set of BE SFs
- $T_f$  : time frame
- $Gr_1^i$  : the amount of bandwidth granted to connection  $i$  during the 1<sup>st</sup> round of the scheduling process.
- $Gr_2^i$  : the amount of bandwidth granted to connection  $i$  during the 2<sup>nd</sup> round of the scheduling process.
- $Gr^i$  : the amount of bandwidth granted to connection  $i$  during the whole grant interval  $I_{gr}^i$ .
- $R_{min}^i$  : The minimum reserved traffic rate for connection  $i$
- $R_{max}^i$  : the maximum sustained traffic rate for connection  $i$
- $L_{max}^i$  : the maximum tolerable latency for connection  $i$
- $I_{gr}^i$  : the grant interval for connection  $i$
- $N_q^i$  : the number of packets in connection  $i$  queue
- $S_q^i$  : the size of connection  $i$  queue in bytes
- $t_{cur}$  : current time
- $t_{lgr}^i$  : time when connection  $i$  got its last grant

The connections participating to the first round of the scheduling process are considered in a strict priority order: UGS, ertPS, rtPS, and nrtPS. Only the amount of data conforming to the minimum rate i.e. equivalent to the number of tokens in the first bucket is scheduled after checking that there is enough bandwidth to carry the corresponding payload and overhead. Note that at the BS side, since different flows may use different MCSs, a translation of  $Gr_1^i$  in terms of time slots/OFDM symbols is needed to evaluate the remaining bandwidth  $BW_r$  (also considered in time slots in this case) (c.f. line 10 of Algorithm 9 and line 9 of Algorithm 11). It is worth mentioning that in this chapter, we consider a DL/UL ratio of 1:1 which is one of the typical ratios recommended by the WiMAX Forum; unlike in Chapter 5 where the DL/UL were dynamically adjusted to the traffic characteristics.

After the first round of the scheduling process, a second round is triggered by the possible availability of extra bandwidth (remaining from the first phase). The objective of this second round is to share the remaining resources among the different connections. In this second round,

---

**Algorithm 9:** BS DL Scheduler: 1st round**Return:**  $W$  the sum of connections weights to be used in the 2nd round

```

1 Begin
2    $W \leftarrow 0$ 
3   for ( $i = 0; i < 5; i++$ ) do
4     for ( $j = 0; j < N_{SF}^i; j++$ ) do
5        $Gr_1^j \leftarrow 0$ 
6        $w^j \leftarrow 0$ 
7       if ( $t_{cur} - t_{lgr}^j \geq I_{gr}^j$ ) then
8          $tmp\_Gr_1^j \leftarrow \min(S_q^j,$ 
9            $R_{min}^j \times I_{gr}^j - Gr^j)$ 
10         $Gr_1^j \leftarrow \text{ovhd\_avail}(tmp\_Gr_1^j, MCS(j))$ 
11         $BW_r \leftarrow BW_r - Gr_1^j$ 
12         $t_{lgr}^j \leftarrow t_{cur}$ 
13         $w^j \leftarrow \min(S_q^j,$ 
14           $R_{max}^j \times I_{gr}^j - Gr^j) - Gr_1^j$ 
15         $Gr^j \leftarrow 0$ 
16         $W \leftarrow W + w^j$ 
17         $W \leftarrow W + \min(S_q^j, R_{max}^j \times I_{gr}^j - Gr^j)$ 
18 return  $W$ 

```

the bandwidth allocation process is performed according to a simple weighted fair queuing strategy. The weight of each connection corresponds to the content of its queue while not exceeding the boundaries set by its two token buckets. After  $Gr_2^i$  is decided, an amount of tokens—corresponding to the payload scheduled in the 2nd round—is removed from the first and then from the second bucket.

In this second phase of the scheduling process, the BE connections are given, proportionally, as much chance as other types of service flows to compete for available resources which could prevent them from starvation. The remaining needs of each non-UGS connection, i.e. the difference between the queue size and the allocated grants are then translated into bandwidth requests. The details of the proposed algorithm are provided for the BS (in DL) and the SS in Algorithm 12 and Algorithm 13, respectively.

Figure 6.2 illustrates the three possible configurations of the token buckets at the end of the grant interval for a given connection  $i$ , after performing the two scheduling rounds. Note that during the whole interval, the buckets are not refilled. In the first case, both buckets are empty which means that the connection has been scheduled at its maximum sustained traffic rate  $R_{max}^i$ . When only the first bucket is empty, this means that the connection has been scheduled at a rate  $R^i$ ;  $R_{min}^i \leq R^i < R_{max}^i$ . In other words, the scheduler has managed to meet at least the minimum requirements of the connection in terms of delay and throughput. The third case, shown in Figure 6.2, corresponds to the case where the first bucket is not completely empty i.e.  $R^i < R_{min}^i$ . This means that the available bandwidth was not enough to cover the needs of the connections participating to the 1st round of the scheduling process. In the two first cases, the two buckets associated to the considered connection are refilled with tokens and the grant interval is reset. In the last case however, the same buckets are maintained. Moreover, to reach  $R_{min}^i$ , the connections

**Algorithm 10: SS Scheduler: 1st round****Return:** W the sum of connections weights to be used in the 2nd round

```

1 Begin
2    $W \leftarrow 0$ 
3   for ( $i = 0; i < 5; i++$ ) do
4     for ( $j = 0; j < N_{SF}^i; j++$ ) do
5        $Gr_1^j \leftarrow 0$ 
6        $w^j \leftarrow 0$ 
7       if ( $t_{cur} - t_{lgr}^j \geq I_{gr}^j$ ) then
8          $tmp\_Gr_1^j \leftarrow \min(S_q^j,$ 
9            $R_{min}^j \times I_{gr}^j - Gr^j)$ 
10         $Gr_1^j \leftarrow \text{ovhd\_avail}(tmp\_Gr_1^j)$ 
11         $t_{lgr}^j \leftarrow t_{cur}$ 
12         $w^j \leftarrow \min(S_q^j,$ 
13           $R_{max}^j \times I_{gr}^j - Gr^j) - Gr_1^j$ 
14         $Gr^j \leftarrow 0$ 
15         $W \leftarrow W + w^j$ 
16      else if ( $(i \in R \text{ or } i \in N)$ 
17         $\text{and } (t_{cur} - t_{lgr}^j + T_f \geq I_{gr}^j))$ ) then
18        if ( $\text{unicast\_BR\_Opp} \geq 1$ ) then
19           $\text{send\_standalone\_BR}$ 
20        else if ( $BWr \geq 6$ ) then
21           $\text{/* bandwidth stealing */}$ 
22           $\text{send\_standalone\_BR}$ 
23        else if ( $N_{SF}^0 \geq 1$ ) then
24           $PM\_bit \leftarrow 1$ 
25         $W \leftarrow W + \min(S_q^j, R_{max}^j \times I_{gr}^j - Gr^j)$ 
26  return W

```

**Algorithm 11: BS UL Scheduler: 1st round****Return:**  $W$  the sum of connections weights to be used in the 2nd round

```

1 Begin
2    $W \leftarrow 0$ 
3   for ( $i = 0; i < 5; i++$ ) do
4     for ( $j = 0; j < N_{SF}^i; j++$ ) do
5        $Gr_1^j \leftarrow 0$ 
6       if ( $t_{cur} - t_{lgr}^j \geq I_{gr}^j$ ) then
7          $tmp\_Gr_1^j \leftarrow \min(Req^j,$ 
8            $R_{min}^j \times I_{gr}^j) - Gr^j$ 
9          $Gr_1^j \leftarrow ovhd\_avail(tmp\_Gr_1^j)$ 
10         $BW_r \leftarrow BW_r - Gr_1^j$ 
11         $t_{lgr}^j \leftarrow t_{cur}$ 
12         $w^j \leftarrow \min(Req_q^j,$ 
13           $R_{max}^j \times I_{gr}^j) - Gr^j - Gr_1^j$ 
14         $Gr^j \leftarrow 0$ 
15         $W \leftarrow W + w^j$ 
16        else if ( $(i \in R \text{ or } i \in N)$ 
17          and ( $t_{cur} - t_{lgr}^j + T_f \geq I_{gr}^j$ )
18          and ( $(N_{SF}^0 == 0)$ 
19            or ( $(N_{SF}^0 > 0 \text{ and } PM == 1)$ ))
20          then
21           $Unicast\_Poll$ 
22           $W \leftarrow W + \min(Req^j, R_{max}^j \times I_{gr}^j - Gr^j)$ 
23   return  $W$ 

```

**Algorithm 12: BS DL Scheduler: 2nd round**

```

1 Begin
2    $W \leftarrow 0$ 
3   for ( $i = 0; i < 5; i++$ ) do
4     for ( $j = 0; j < N_{SF}^i; j++$ ) do
5        $tmp\_Gr_2^j \leftarrow \frac{w^j}{W} \times BW_r$ 
6        $Gr_2^j \leftarrow ovhd\_avail(tmp\_Gr_2^j)$ 
7        $BW_r \leftarrow BW_r - Gr_2^j$ 
8        $Gr^j \leftarrow Gr^j + Gr_2^j$ 

```

**Algorithm 13:** SS Scheduler: 2nd round

---

```

1 Begin
2    $W \leftarrow 0$ 
3   for ( $i = 0; i < 5; i++$ ) do
4     for ( $j = 0; j < N_{SF}^i; j++$ ) do
5        $Gr_2^j \leftarrow 0$ 
6       if ( $w^j > 0$ ) then
7          $tmp\_Gr_2^j \leftarrow \frac{w^j}{W} \times BW_r$ 
8          $Gr_2^j \leftarrow ovhd\_avail(tmp\_Gr_2^j)$ 
9          $BW_r \leftarrow BW_r - Gr_2^j$ 
10         $Gr^j \leftarrow Gr^j + Gr_2^j$ 
11       if ( $Gr_2^j > 0$  and  $S_q^j > 0$ ) then
12         if ( $BW_r > 2$ ) then
13            $Piggyback\_BR$ 
14         else if ( $Contention\_BR\_Opp$ ) then
15            $send\_standalone\_BR$ 

```

---

needs more bandwidth than what is reflected by the content of the first bucket. Therefore, at the beginning of the following frame  $T_f \times R_{min}^i$  tokens from the second bucket are marked indicating that the threshold for the 1st round is not only set by the content of the first bucket but also with the marked tokens from the second one. The connection participates to the first round of the scheduling process as many times as needed, during the following time frames, till all the tokens of the first bucket and those marked in the second bucket are removed. It is only at that time that the two buckets associated to this connection are refilled and the grant interval is reset. This last case entails some latency for the considered flow. Nevertheless by shifting the corresponding grant interval, we decrease the chances that the same thing happens again (two or more heavy bursts coincide in the same time frame) especially if the burstiness occurs periodically.

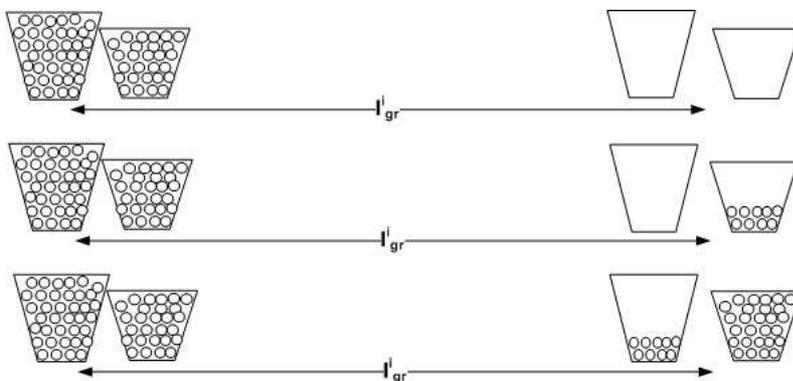


Figure 6.2: A Dual-Bucket Shaping Mechanism

### 6.2.1 Bandwidth request and grant strategy

As reported in Algorithms 11, 10, and 13, the bandwidth request and grant strategy we adopt in this chapter is the following:

- The BS polls individually each rtPS or nrtPS connection whose grant interval is expiring within one time frame. These unicast bandwidth request opportunities are allocated (as long as possible) in the contention interval. If all the slots of the contention interval have been used, the BS allocates these bandwidth request opportunities in the data grant interval which would allow the MS to perform bandwidth stealing if required.
- The rest of bandwidth request contention interval slots, if any, are addressed to a group of MSs. This group consists of all the MSs which have at least one nrtPS or BE connection, which are the only types of services allowed to use group and/or broadcast bandwidth request opportunities.
- Note that, as mentioned in the standard, the periodic unicast polling concerns only MSs which do not have any active UGS connection, unless the PM bit is set. Moreover, because setting the PM bit to 1 does not tell how many rtPS/nrtPS connections need to be polled, we have chosen to interpret it as a request from the MS to have unicast bandwidth request opportunities for each of its polling service connections whose grant interval is expiring within two time frames.
- As we have seen in Section 3.1.3, piggybacking is possible only when packets from the same connection, that is requiring more bandwidth, are transmitted. Therefore, the MS uses this technique in the second round of the scheduling process when a connection is scheduled but still has packets in the queues (c.f Algorithm 13 lines 11-13).

The choice of one or another of the available techniques is motivated by the concerned service type and by the overhead entailed by the use of that technique.

## 6.3 Performance Analysis

To evaluate the performance of mCoSS, we have implemented the corresponding set of algorithms under QualNet 4.5 [31] which is the commercialized version of GloMoSim. mCoSS has been compared to SP and to a variant of the WFQ discipline. In this section we first give an overview of the features supported by the WiMAX simulation model proposed by QualNet. Then, we define the scenarios and simulation settings considered in the performance analysis before reporting and commenting the obtained results.

### 6.3.1 A WiMAX simulation model under QualNet

QualNet 4.5 provides the Advanced Wireless Model Library which addresses both fixed and mobile WiMAX systems. The proposed simulation model is dedicated to OFDMA-based PMP networks operating in TDD mode. It supports the five service types UGS, ertPS, rtPS, nrtPS and BE and several types of bandwidth request mechanisms (polling-based, contention-based and CDMA-based). Most of the IEEE 802.16 management messages (DCD, UCD, UL-MAP, DL-MAP, DSx, etc.) are implemented and several features like the AMC, fragmentation, and packing are supported. Nevertheless, some bugs in the fragmentation mechanism (leak in the queues) have been noticed. We have fixed this bug by correcting the way the queue size is updated when a fragment

or a whole packet is removed from the queue. Moreover, only CBR and VBR generators have been considered when mapping the QoS parameters from application to MAC level. We have extended this capability to Super-Application traffic generator which provides more flexibility in the flow configuration. The model provides also a basic admission control mechanism and a scheduling policy based on a variant of the WFQ strategy, which is different from the one we use in mCoSS. The WFQ variant implemented in QualNet calculates and assigns a finish time to each packet. In this calculation, WFQ uses the bit rate of the link, the number of queues, and the size of each packet in each of the queues. The WFQ scheduler then transmits the packet with the earliest finish time among all the queued packets. Thus, each time a packet is dequeued, the WFQ scheduler recomputes the finish time assigned to each packet which entails a high computational complexity and limits the scalability of the proposed approach.

### 6.3.2 Performance evaluation

Channel Frequency	3.5 GHz
Channel bandwidth	10 MHz
FFT size	2048
Cyclic prefix gain	8
Propagation pathloss model	Two-ray
BS antenna Tx power	33 dBm (= 2 W)
BS antenna height	32 m
BS antenna gain	15 dBi
MS antenna Tx power	23 dBm (= 200 mW)
MS antenna height	1.5 m
MS antenna gain	-1 dBi
Type of antenna	omnidirectional
Frame duration	10 ms
DL subframe duration	5 ms

Table 6.1: Simulation settings

In this section, we consider the parameters settings reported in Table 6.1. As mentioned before, we consider a DL/UL ratio of 1:1 from a total frame size of 10 ms. A simple two-ray pathloss propagation model has been used and no shadowing or fading has been considered to offer a "simple" environment for the comparison of the different algorithms.

In the following scenarios, we consider an audio stream of 30 mns configured as an UL rtPS connection. The audio frame size is set to 1600 bytes and the number of frames per second follows a uniform distribution between 10 and 25 fps (frame/second). The QoS parameters of the considered stream are configured as follows:  $R_{min}^i = 128$  kbps,  $R_{max}^i = 320$  kbps and  $I_{gr}^i = 100$  ms.

#### 6.3.2.1 Scenario 1: mCoSS shaping capability

In this scenario, we propose to test the shaping capability of our multi-Constraints Scheduling Strategy (mCoSS). Therefore, we place two MSs at the same distance from the BS and we configure an audio stream for each MS as mentioned before:  $R_{min}^i = 128$  kbps,  $R_{max}^i = 320$  kbps and  $I_{gr}^i = 100$  ms. While MS1 respects these boundaries, MS2 transmits the audio stream at a much higher rate varying from 640 kbps to 1.28 Mbps. More than 30 experiments have been run to

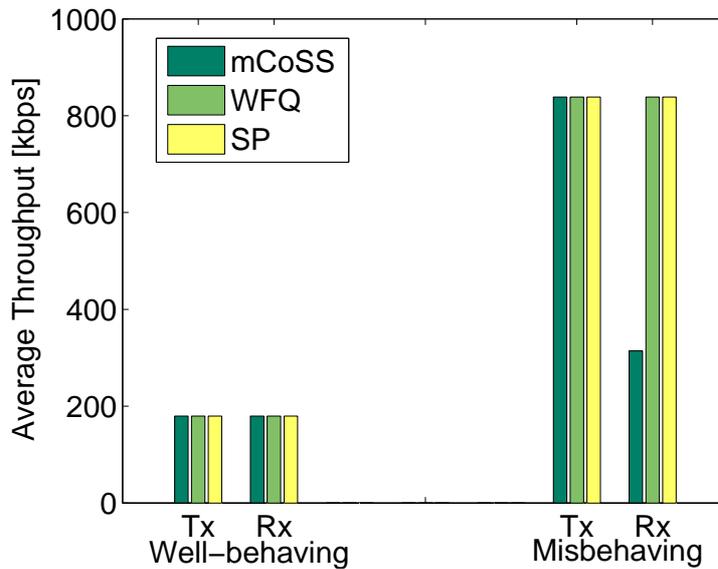


Figure 6.3: mCoSS Shaping Capability

	MS1 Well-behaving	MS2 Misbehaving
mCoSS	0.255	13.6
WFQ	0.57	0.53
SP	0.57	0.53

Table 6.2: mCoSS Shaping Capability: E2E Delay (sec)

validate the shaping capability of our algorithm and to compare it to the WFQ and SP algorithms implemented in QualNet. Figure 6.3 plots the transmission (Tx) and reception (Rx) rates of both the misbehaving and the well-behaving traffics for the three algorithms: mCoSS, WFQ, and SP. The Tx rate represents the rate at which the application is generated at the MS while the Rx rate is the reception rate at the BS. We can see from Figure 6.3 that for the well-behaving traffic sent by MS1, the three algorithms have almost equal performance in terms of throughput. For the misbehaving traffic however, both SP and WFQ allow it to reach more than 800 kbps while mCoSS forces the traffic to stay within the set boundaries: the reception rate at the BS does not exceed 315 kbps. Tables 6.2 and 6.2 report the obtained E2E delay and jitter for both traffics using the different scheduling algorithms. As a consequence of the policing/shaping policy adopted by mCoSS, the misbehaving traffic generated by MS2 is penalized (in comparison to SP and WFQ) in terms of E2E delay since packets exceeding  $R_{max}^i$  are delayed and possibly dropped if their number exceed the buffers capacity. On the other hand, the E2E delay of well-behaving traffic is halved compared to WFQ and SP. With both WFQ and SP the two traffics experience comparable E2E delays; the misbehaving traffic gets even a shorter average jitter than the well-behaving traffic.

From the obtained results, we can see that mCoSS is capable of forcing a traffic to stay within the allowed thresholds and of isolating a well-behaving traffic from a misbehaving one. The absence of shaping at WFQ and SP has affected the performance of the first traffic and could even have a much worse effect if the second traffic had been generated at a rate that overload the whole network.

	MS1 Well-behaving	MS2 Misbehaving
mCoSS	22	80
WFQ	69	27.7
SP	69	27.7

Table 6.3: mCoSS Shaping Capability: Jitter (ms)

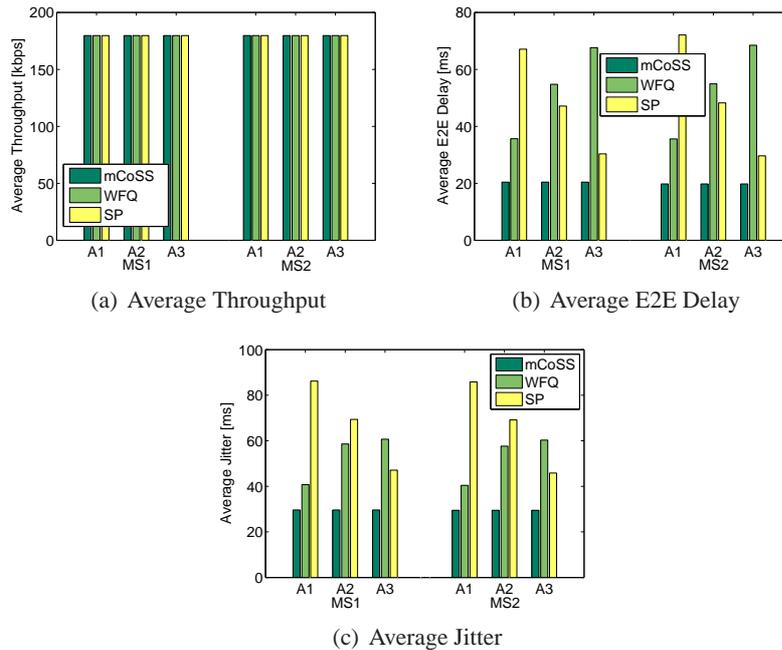


Figure 6.4: 2 MSs with 3 Audio streams each

### 6.3.2.2 Scenario 2: fairness and QoS degree of satisfaction

In this second scenario, we consider the same MSs having each three audio streams with the same configuration. Through this scenario, we aim at evaluating, in same channel and traffic conditions, the performance of our scheduling algorithm in terms of inter-MSs and inter-SFs fairness and to compare the QoS degree of satisfaction of the six connections using the three algorithms. Figure 6.4(a) plots the obtained average throughput of the 1st, 2nd, and 3rd audio streams (A1, A2, and A3) of MS1 and MS2. As long as the throughput is concerned, the three algorithms offer the same level of performance. The average end-to-end (E2E) delay and jitter, however, experience a less stable behavior from one algorithm to another as we can see from Figures 6.4(b) and 6.4(c), respectively. With WFQ, the E2E delay varies from 35 to 67 ms from one service flow/MS to another. The same behavior is noticed for SP for which the E2E average delays vary from 30 to 72 ms. mCoSS on the other hand provides lower and much more stable results for the six flows for both E2E delay (around 20 ms) and jitter (less than 30 ms).

Considering throughput, delay and jitter, mCoSS, in comparison with SP and WFQ, provides the best and most stable performance among SFs which results in a better inter-SFs and inter-MSs fairness.

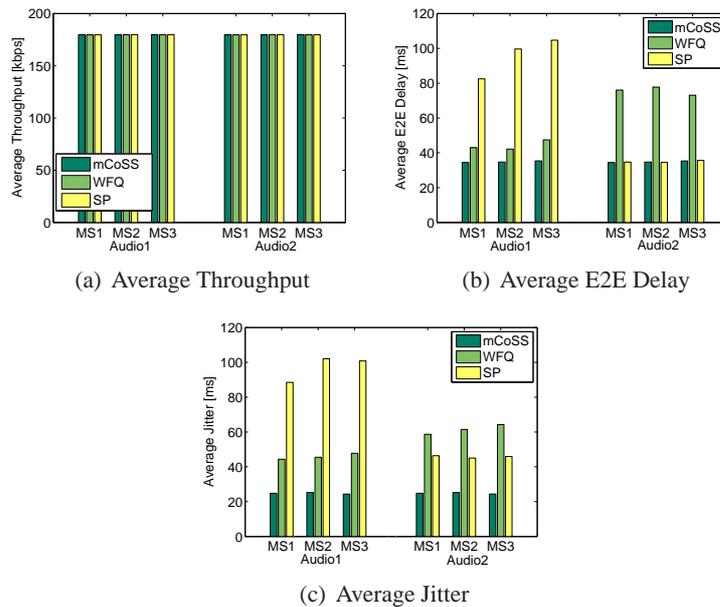


Figure 6.5: 3 MSs with 2 Audio streams each

### 6.3.2.3 Scenario 3: AMC support

Through this last scenario, we aim at validating the capability of mCoSS to adapt the allocated bandwidth to the channel conditions of the MS; a capability that is already supported in QualNet implementation of WFQ and SP. Since the objective is to test the AMC capability of mCoSS, we consider 3 MSs placed at three different positions from the BS: at 1km, 2 km, and 3 km away from the BS. These three distances correspond to an three SNR levels matching UIUC 1 (QPSK 1/2), UIUC 4 (16-QAM 3/4) and UIUC 7 (64-QAM 3/4). We configure two audio streams at each MS with the same settings previously specified. As we can see from Figure 6.5(a), like for the previous scenario, the three algorithms have almost equivalent performance for the throughput. However, the difference in E2E delay (plotted in Figure 6.5(b)) between Audio 1 and Audio 2, with SP, is more noticeable than in the second scenario. Indeed it varies for MS3 for example from 35 ms to more than 100 ms which exceeds the maximum latency of the service flow. The same behavior is observed for jitter in Figure 6.5(c).

For mCoSS, the increase of the number of MSs and the use of different MCSs had almost no effect on the performance of the algorithm. The same stability of results is observed in this scenario which confirms the fairness of the algorithm and its scalability at low level at least.

The performance evaluation of mCoSS presented in this chapter is by no means comprehensive, yet it shows and validates some of the key features supported by the proposed strategy. More simulation scenarios though—involving more service types—need to be considered.

## 6.4 Conclusion

Most of the hierarchical scheduling strategies proposed in the literature and described in Chapter 4 (such as [13, 10, 9]) propose a specific queuing discipline for each scheduling service type, which increases significantly the complexity of the proposed scheduling policy. Unlike those approaches, the multi-Constraints Scheduling Strategy (mCoSS) proposed in this chapter is designed to be applicable to all service types. Based on a modified dual-bucket traffic shaping mechanism used for

all the scheduling service types, mCoSS allies the genericity of the approach to the specificity of the configuration since the dual-bucket mechanism is configured on a per-flow basis.

This shaping mechanism is combined with a two-rounds scheduling strategy which reflects (i) at the first round, the minimum data rates and latency requirements the BS or MS is committed to provide and (ii) at the second round, the efficiency and fairness of the resources management since the remaining bandwidth is shared in this round using a simple WFQ strategy; the allocations should nevertheless remain within the thresholds set by the dual-bucket shaping mechanism. The bandwidth request and grant mechanism adopted in mCoSS is designed to make a tradeoff between increasing the accuracy of the bandwidth needs perception at the BS and decreasing the overhead associated to frequent unicast polling. Indeed the proposed strategy alternates between bandwidth stealing, piggybacking, unicast, broadcast and group polling, and the use of PM bit based on the considered scheduling service type and the available resources. The proposed mCoSS has been implemented under QualNet 4.5 simulator and compared to SP and to a variant of WFQ discipline. The preliminary results reported in this chapter validate and confirm the shaping fairness and AMC support capability of the proposed mCoSS. They also show that, compared to SP and WFQ, mCoSS provides better and more stable end-to-end-delay and jitter performances. More simulations need though to be carried out to check and validate other aspects of the proposed scheduling strategy.

---

## Chapter 7

# Mobile WiMAX: a V2I Communications Medium

Intelligent Transportation Systems (ITS) have been under development since the 80's as part of a global strategy for solving many of our modern life transportation problems. These systems enable people to reach their destinations in a safe, efficient, and comfortable way. In order to reach that goal, several radio access technologies (RAT) such as UMTS, WiFi, WiMAX and 5.9 GHz have been proposed for next generation ITS.

In addition to the 5.9 GHz, which is dedicated to vehicular ad hoc networks networks, mobile WiMAX is expected to play a major role in ITS since it is the only mobile broadband technology currently in use.

Yet, the coexistence between 5.9 GHz technology, mobile WiMAX, and other technologies in the vehicles raises the challenge of choosing the most appropriate RAT. In order to address this problem and define optimal rules for the communication technology selection, comparisons on the network performance have to be done.

In this chapter, we compare mobile WiMAX (based on IEEE 802.16e standard) and 5.9 GHz technology (based on the upcoming IEEE 802.11p standard). We investigate, through simulation, the potential and limitations of both technologies as a communication media for vehicle-to-infrastructure (V2I) communications. The performance of the two systems is evaluated for different vehicle speeds, traffic data rates, and network deployments

The remainder of this chapter is structured as follows. Section 7.1 presents the categories of applications targeted by the intelligent transportation systems. Section 7.2 provides an overview of IEEE 802.11p, and summarizes the main characteristics of mobile WiMAX and 5.9 GHz technologies and compares them based on several criteria. In Section 7.3, we first define our simulation environment and settings and then analyze the results of the performance evaluation study we have performed. Section 7.4 concludes the chapter by outlining the main obtained results.

### 7.1 ITS applications and architectures

During the last two decades, several initiatives, like COMeSafety [57], and technical groups supported by standardization bodies, such as the IEEE 802.11p task group [58], the ISO TC204 Working Group 16 [59] and the ETSI ITS Technical Committee [60] have been created to solve many of our society transportation problems. From that perspective, three main categories of applications have been targeted: (i) road safety applications, (ii) traffic efficiency applications, and (iii) value-added applications.

---

Table 7.1: ITS Applications categories: examples and requirements.

Application category	Latency tolerance	Range	Example (delay requirements)
Road safety	Low latency	Local range	Pre-crash sensing/warning (50 ms) Collision risk warning (100 ms)
Traffic efficiency	Some latency is acceptable	Medium range	Traffic information - Recommended itinerary (500 ms)
Value-added services	Long latency is accepted	Medium range	Map download update - Point of interest notification (500 ms)

- Road safety applications: the primary goal of this set of applications is to reduce road fatalities by assisting and warning the driver about the potential risks. This category covers applications like pre-crash sensing and collision risk warning.
- Traffic efficiency applications: this category is intended to relieve traffic congestion by helping to monitor the traffic flow and by providing alternative itineraries to drivers. These applications make the transportation systems not only more efficient but also more environmentally friendly by optimizing routes and decreasing gas emissions.
- Value-added applications: they include on-demand services related to infotainment, comfort or vehicle management. They can be provided either free of charge or for a fee - which could help to finance the deployment of such networks. Also, by notifying a point of interest (e.g. parking lot, restaurant, etc.), some of these applications may help to save time and thus to reduce fuel consumption.

In Table 7.1 we can see that the groups of services presented above have different requirements, in terms of range, delay, and throughput. Indeed, they cover a wide range of applications that vary from "locally" sending a small and urgent message (e.g., in order to alert a driver about an imminent crash) to updating a map on the on-board device by downloading a big file from a remote server. Considering the conflicting requirements of the applications, several ITS architectures have been proposed by vehicular communications initiatives and standardization bodies. In particular, most of them agree on the necessity of having a variety of communication media. The two architectures, presented in Figures 7.1(a) and 7.1(b), are proposed by the European Telecommunications Standards Institute (ETSI) [57], and ISO TC204 Working Group 16 [59], respectively.

The possibility of having different communication technologies for vehicular communication yields to the necessity to understand which is the most suitable in every specific context. Indeed, since in the near future vehicles will be equipped with different access technologies, knowing the capabilities and limitations of these technologies, and knowing their availability are very important factors to make radio access technology (RAT) selection and decide whether a vertical handover should be performed to achieve an always best connected communication.

Recently, standardization bodies have given mandate to technical groups to define the application requirements for ITS applications. Moreover, business models will be developed to include the cost and benefit for the user of using a certain technology with respect to another. The last piece needed is the performance analysis of the different access technologies.

Among the communication technologies, in this chapter we propose to compare two of the most promising ones: mobile WiMAX (based on IEEE 802.16e standard [2]) and the 5.9 GHz technology based on the upcoming IEEE 802.11p standard.

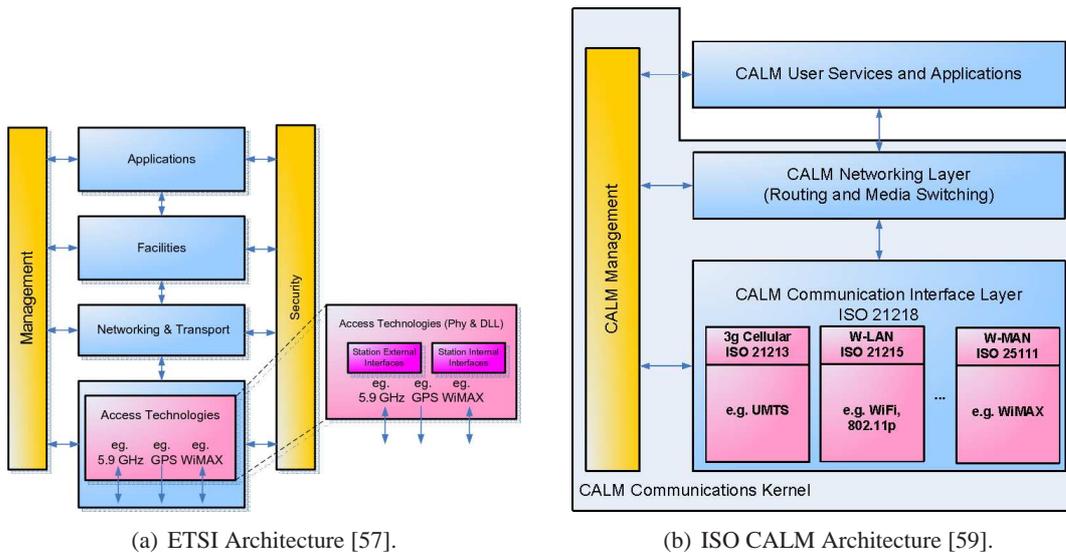


Figure 7.1: ITS station reference architectures.

IEEE 802.11p-based technology [58] has been developed for the specific context of vehicular networks. It is expected to be particularly suitable for medium range and delay-sensitive road safety applications. Mobile WiMAX, on the other hand, offers a promising alternative because of its potential to offer medium to long range connectivity, full support of mobility, and high data rates with moderate delay.

Based on these characteristics, the two technologies seems intrinsically complementary in terms of range, data rates and delay. Nonetheless, to the best of our knowledge, this is the first time that the performance of the two technologies are compared through simulation. Our objective is to study the feasibility of both technologies as communication media for vehicular networks by evaluating their performances in the same simulation environment.

## 7.2 IEEE 802.11p vs. IEEE 802.16e

IEEE 802.11p is an ongoing 802.11 amendment [58] that is aimed at standardizing a set of extensions for 802.11 in order to adapt it to the V2X (vehicle-to-infrastructure V2I and vehicle-to-vehicle V2V) environment.

From that perspective, many phases of the basic 802.11 communication protocol at MAC layer have been eliminated or shortened. Indeed, unlike 802.11, 802.11p allows stations to communicate in OCB mode i.e. outside the context of a basic service set (BSS), thus avoiding the latency caused by the association phase. Moreover, there is no need to scan the channel since the OCB communication occurs in a frequency band dedicated to ITS use<sup>1</sup>. Also, when exchanging frames in OCB mode, the MAC layer authentication services are not used. Yet, it is still possible to have secured communications provided by mechanisms outside the MAC layer.

At physical layer, the amendment concerns mainly the spectrum allocation. Vehicular communications are performed in the 5 GHz range, where one channel is dedicated to control and the others to ITS services. Figure 7.2 illustrates in particular the European profile for the channel allocation. According to this profile, the control channel (G5CC) is used for road safety and traffic efficiency applications. It may also be used to announce ITS services operated on the service

<sup>1</sup>A license might be needed for these bands, depending on the regulatory domain.

channels (G5SC1 to G5SC5). The service channels G5SC1 and G5SC2 are used for ITS road safety and traffic efficiency applications while the others (G5SC3, G5SC4 and G5SC5) are dedicated to other ITS user applications. In order to reduce the effects of Doppler spread, the use of 10 MHz channels has been adopted instead of the usual 20 MHz used by 802.11a. Consequently, all OFDM timing parameters are doubled (e.g. the guard interval, the OFDM symbol duration, etc.) and the data rates are halved (vary from 3 to 27 Mbps instead of 6 to 54 Mbps). Moreover, the European profile requires that ITS stations are able to simultaneously receive on both the control and one service channel. Therefore, two transceivers are needed. In this work, we considered the standard profile of the physical and MAC layers recently proposed by ETSI [32].

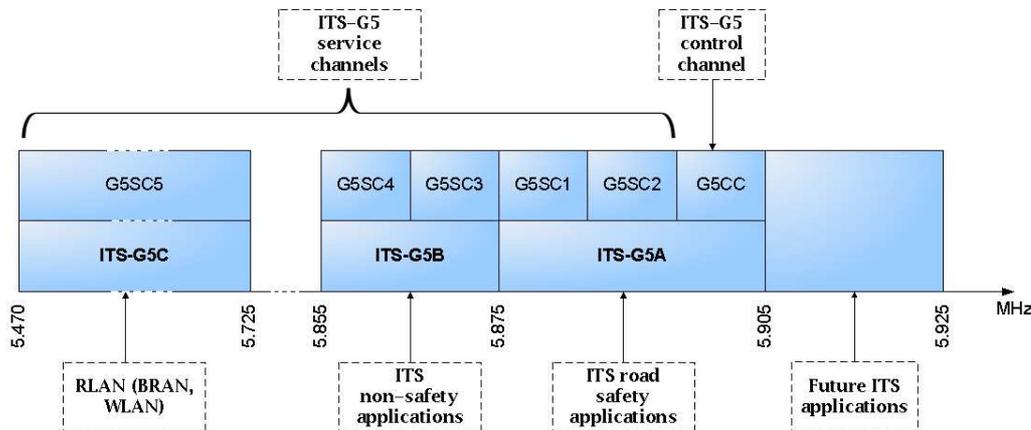


Figure 7.2: European channel allocation [32].

Table 7.2 summarizes the characteristics of both technologies based on several criteria including the frequency spectrum in use, the medium access technique, and the support of security and QoS.

## 7.3 Performance evaluation

### 7.3.1 Simulation environment and settings

For our simulations, we have used the network simulator QualNet 4.5 [31] which is the commercialized version of GloMoSim. The Advanced Wireless Library proposed by QualNet integrates a simulation model for mobile WiMAX with the support of several features such as PHY OFDMA, PMP and TDD modes, AMC capability, QoS scheduling services, etc. Nevertheless, the simulator does not include an 802.11p model. Therefore, we have first implemented the necessary changes (as reported in Section 7.2) to existing 802.11a PHY and 802.11e MAC models in order to adapt them to 802.11p specifications. Note that we have adapted the power of the transmitter and the minimum sensitivity of the receiver to what has been specified in [32].

To evaluate and compare the performance of both mobile WiMAX and 802.11p technologies in V2I context we have considered a highway scenario. Our study is divided in three parts. During the first part we measure the connectivity of the two technologies in order to determine the radio range between a vehicle and a 802.11p road side unit (RSU), or a WiMAX base station (BS). In the second part, we compare the communication performance of the two technologies on a highway segment the length of which corresponds to the coverage of one BS varying the speed of the vehicle. After analyzing the performance of WiMAX, the performance of 802.11p is investigated

Table 7.2: 802.11p vs 802.16e

	802.11p	802.16e
<b>Standardization</b>	Draft [58]	Standard [2]
<b>Frequency/ License</b>	5.470-5.925 GHz free but licensed “License by rule”	10-66 GHz licensed below 11 GHz: (2.3, 2.5, 3.5, 5.8, etc.) both licensed and license-exempt
<b>Channel bandwidth</b>	10 MHz	Depends on the Phy profile (3.5, 5, 7.5, 10 MHz, etc.)
<b>QoS support</b>	4 classes of QoS (EDCA extension) AC_VO, AC_VI, AC_BK, AC_BE	5 classes of QoS: UGS, ertPS, rtps, nrtPS, BE.
<b>Security support</b>	No Authentication prior to data exchange Instead, each packet is used for authentication by certificate based digital signatures	data encapsulation protocol with a set of cryptographic suites and PKM protocol to synchronize keying data between BSs and MSs
<b>Media access technique</b>	CSMA/CA No scanning, no association	TDMA, FDD or TDD
<b>Usage</b>	Network dedicated to vehicles (ITS stations)	Could be used by residences, companies, personal devices, ...
<b>Other supported features</b>		Support of AMC, ARQ, AAS, STC and MIMO

by replacing the single BS with the number of RSUs necessary to cover the same segment. Finally, in the third part, we observed the impact of the traffic datarate and the vehicle speed on the throughput and the delay.

In order to determine the range of an 802.11p RSU and of a WiMAX base station, we have set our simulation parameters as reported in Table 7.3. The pathloss fading model has been set to a two-ray Ricean fading model with a high line-of-sight component which is quite realistic in the highway context (unlike in an urban environment, where this assumption is not valid).

For the evaluation of the range of an 802.11p RSU, we simulated the transmission of periodic beacons (using the control channel at 5.9 GHz for 802.11p communication). Accordingly to the ETSI specifications, the basic beaconing rate is set to 10 Hz and the periodic message (also called CAM, i.e. cooperative awareness message) is 55 bytes long and contains geo-information. The scenario is illustrated in Figure 7.3(a).

In Figures 7.3(c) and 7.3(d), we can observe the delivery ratio as a function of the vehicle distance from the RSU or the BS. Considering a packet delivery ratio greater than 90%, the cell radius coverage of 802.11p and WiMAX are then around 900 meters and 6.5 Km, respectively.

Based on these results, we have set three different network deployments for all the simulation scenarios to be considered. The first deployment corresponds to the case of a highway of 13 km fully covered by one WiMAX base station. The second deployment consists in fully covering the same road link by the equivalent number of 802.11p RSUs (as shown in Figure 7.4(a)). Finally, in order to observe the effect of handover on mobile WiMAX performance too, we have considered a third deployment that considers the area covered by two WiMAX BSs.

Table 7.3: Simulation parameters

	802.11p	802.16e
<b>Frequency</b>	5.87 GHz (G5SC3)	3.5 GHz
<b>Channel bandwidth</b>	10 MHz	10 MHz
<b>RSU Tx power</b>	23 dBm (=200 mW)	33 dBm (=2 W)
<b>RSU antenna height</b>	2.4 m	32 m
<b>RSU antenna gain</b>	3 dBi	15 dBi
<b>MS Tx power</b>	23 dBm (=200 mW)	23 dBm (=200 mW)
<b>MS antenna height</b>	1.5 m	1.5 m
<b>MS antenna gain</b>	0 dBi	-1 dBi
<b>Type of antenna</b>	Omnidirectional	
<b>Pathloss</b>	Two-ray	
<b>Fading model</b>	Ricean	

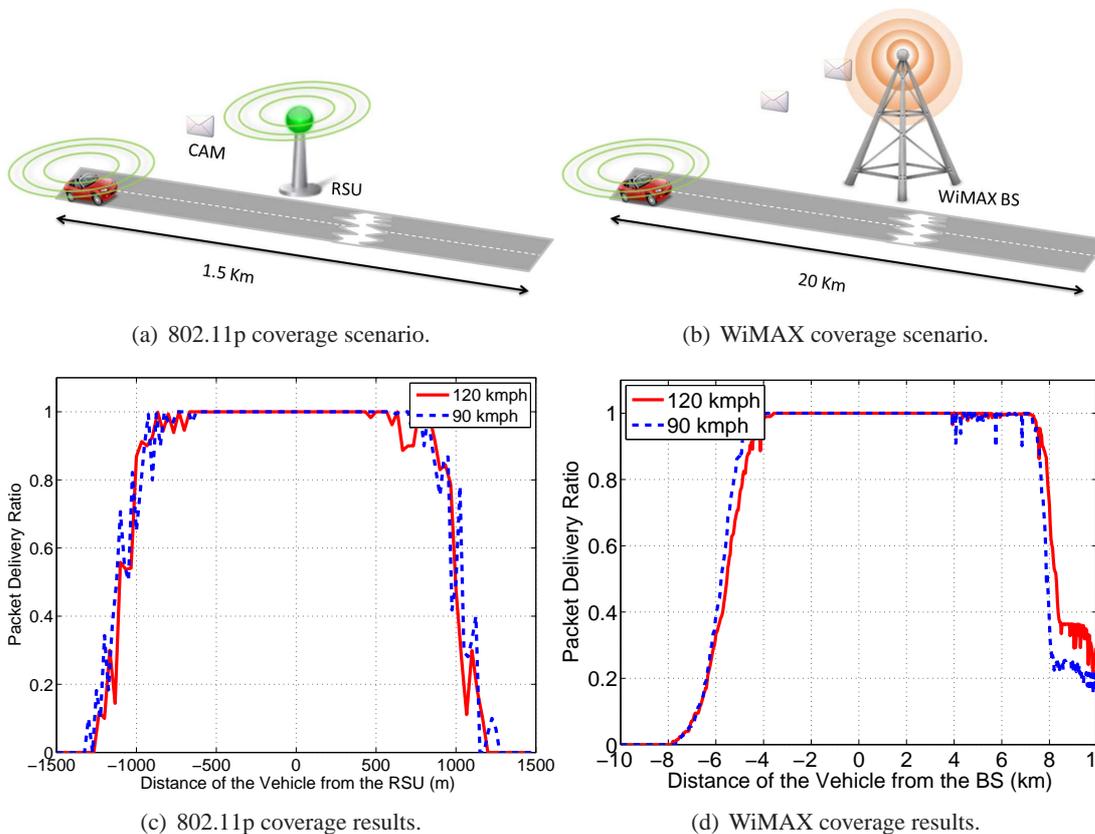


Figure 7.3: Coverage evaluation scenarios.

In all the scenarios, we have considered a source of traffic that is connected to the RSUs/BSs through Ethernet links of 100 Mbps (to avoid any bottleneck outside the considered WiMAX/802.11p V2I network). In the case of 802.11p scenarios, we simulated the transmission of the data over the G5SC3 channel, which is dedicated to non-safety applications.

The effect of increasing the number of vehicles is not considered in this chapter. In fact, even with only one vehicle, by increasing the source data rate, we can analyze the upper limits that can be reached in mobile WiMAX and 802.11p V2I networks in similar conditions.

In order to have realistic movement of the vehicle on the highway, the mobility traces have been generated with SUMO 0.9.8 [61]. In particular, in order to adapt the mobility traces generated by SUMO to QualNet, we have used MOVE (MOBility model generator for VEhicular networks) tool [62].

### 7.3.2 Performance analysis

Using the simulation parameters detailed in Section 7.3.1, we have considered two scenarios.

#### 7.3.2.1 Scenario 1: Study of the impact of the source data rate on 802.11p/802.16e V2I networks performance

In this first scenario, we have set the average speed of the vehicle to 100 kmph, that is a realistic value of vehicles on the highway. We have varied the data rate of a CBR traffic transmitted from the source to the vehicle considering the three configurations of deployed networks. This scenario covers network traffic loads varying from 25 kbps to 20 Mbps. We have evaluated the impact of varying the source data rate on both the throughput (shown in Figure 7.5(a)) and the end-to-end delay (illustrated in Figure 7.5(b)). In the case of 802.11p, we investigated the impact of using RTS/CTS on the transmission performance. In fact, the ETSI standard [32] allows the use of this mechanism for unicast transmissions whose packet size exceeds the *dot11RTSThreshold*. Thus, giving that the packet size is set to 512 bytes, we considered two cases; first the *dot11RTSThreshold* is set to 0 and then to 1000 bytes, which is the default value recommended by ETSI.

All the results presented in this Section are the values averaged over more than 30 runs within a 95%-confidence interval.

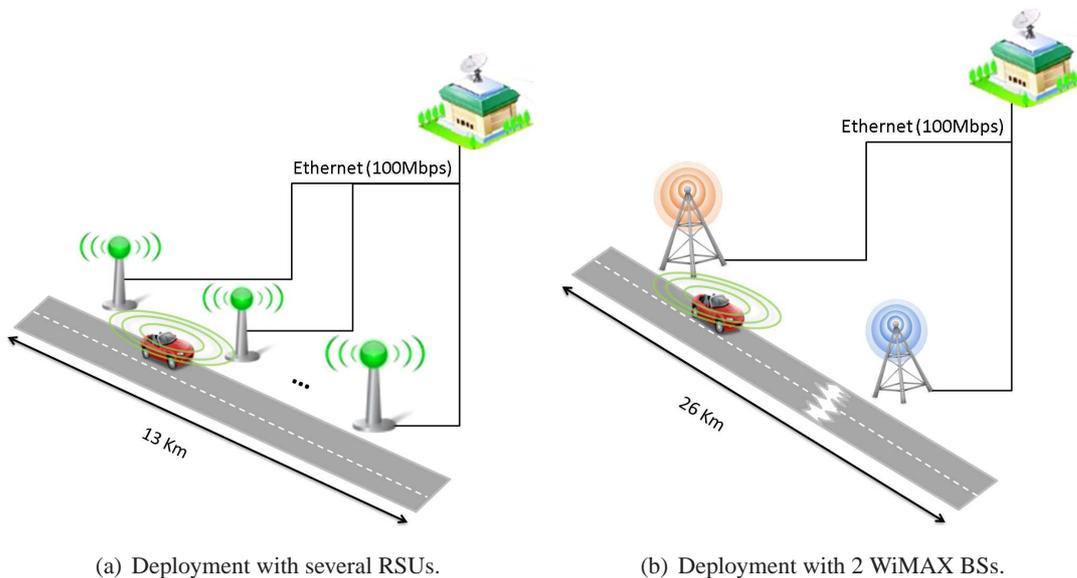


Figure 7.4: Scenarios network deployments.

The obtained results allow us to derive the maximum throughput that could be reached in optimal (1 vehicle) yet realistic conditions (of speed, power, fading, etc). For IEEE 802.11p, the maximum throughput is around 1.2 Mbps while it could exceed 12 and 13 Mbps in 2 BSs and 1 BS deployment scenarios, respectively. As for the average end-to-end (E2E) delay, 802.11p experiences short delays (less than 40 ms) in low traffic conditions. However, when the source

data rate exceeds the maximum that could be reached in 802.11p networks (around 1.2 Mbps), the delay significantly increases, exceeding 200 ms. When using RTS/CTS mechanism the delay further increases. The same behavior (increase of the E2E delay) is observed for WiMAX when the maximum sustainable data rate is reached, though at much lower scale since the average delay does not exceed 60 ms which fulfills even the needs of most emergency applications. However, at very low data rate (e.g. 25 kbps) 802.11p performs better than 802.16e which is convenient for exchanging small and delay-sensitive safety messages.

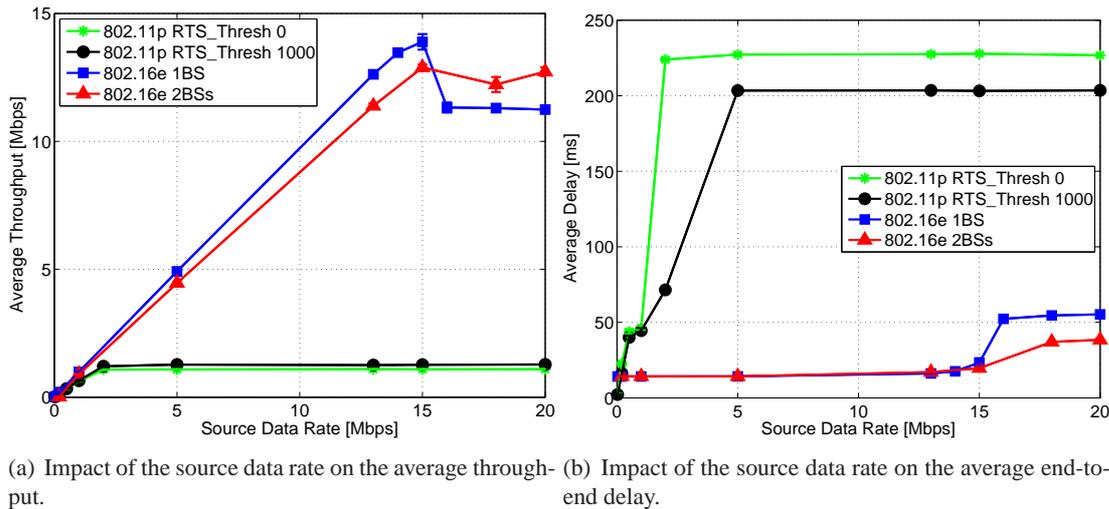


Figure 7.5: Impact of the source data rate on the average performance

### 7.3.2.2 Scenario 2: Study of the impact of the vehicle speed on 802.11p/802.16e V2I networks performance

In this second scenario, we have set the source data rate to 1 Mbps, a value that is slightly below the limit of 1.2 Mbps that we observed in the previous scenario for 802.11p case, but that should maintain a good throughput. We have observed the impact of varying the vehicle speed on the average throughput (plotted in Figure 7.6(a)) and the end-to-end delay (shown in Figure 7.6(b)).

For 802.11p, when the vehicle speed increases, the connectivity time to the 802.11p RSUs decreases which then reduces the amount of data received by the vehicle. Additionally, a fraction of time of this period is required to switch from one RSU to another. On the other hand, in the case of two WiMAX BSs, the handover execution requires a non-negligible time which affects the average throughput that remains lower than that of the scenario with a single BS regardless of the vehicle speed.

The average E2E delays of 802.11p and 802.16e are plotted in Figure 7.6(b)). Remind that in this scenario, the source data rate is set to 1 Mbps, so there is no packet loss due to buffer overflow at the IP or MAC layers. For this reason, the end-to-end delay is the same with one and two WiMAX base stations while in case of 802.11p, the delay slightly increases with the vehicle speed. One important observation that could be derived from this figure is that for both technologies, the E2E delay is lower than 55 ms (less than 15 ms for mobile WiMAX) which fulfills the minimum requirement of most ITS safety applications. As final remark, the use of RTS/CTS mechanism slightly increases the E2E delay and affects the throughput. Nevertheless, the impact of this mechanism should be further investigated in heavy loaded vehicular traffic scenarios where it could prevent collisions and increase the packet delivery ratio but also entail longer delays.

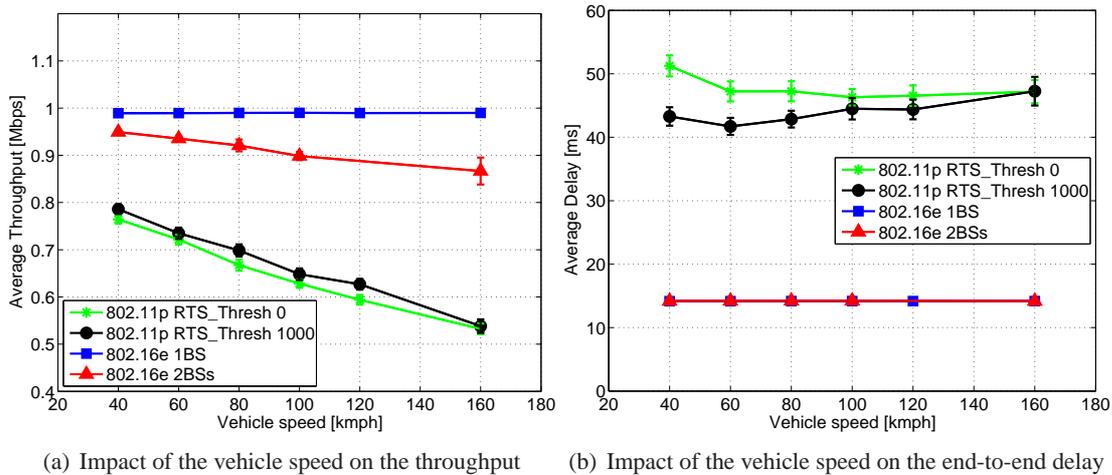


Figure 7.6: Impact of the vehicle speed on the average performance

## 7.4 Conclusion

In this chapter, we studied the potential and limitations of mobile WiMAX as a communication medium for vehicle-to-infrastructure (V2I) communications and more specifically in comparison with 802.11p. We first compared the two technologies based on different criteria. Moreover, we investigated their performance through simulation. The coverage, average throughput, and end-to-end delay were evaluated for different vehicle speeds, traffic data rates, and network deployments.

The simulation results reveal on one side the great competitiveness of mobile WiMAX technology in the context of V2I communications. In particular, this technology, offers, not only a large radio coverage and high data rates, but also reasonable and even very low delays. On the other side, the 802.11p technology is better suited to low traffic loads, where it offers very short latencies even at high vehicle speed.

The obtained results can be considered as a first step for the definition of an efficient common radio resource management (CRRM) module for vehicular networks. They could further be used as pre-defined criteria for radio access technology (RAT) selection for ITS applications. A broad analysis of the performance of the two technologies could be used to develop new algorithms for smart selection of the optimal RAT based on the applications requirements, the channel load, and the user's preferences.



## Conclusion

WiMAX technology, which has emerged as a competitive alternative to wireline broadband access solution, provides QoS support for heterogeneous classes of traffic with different QoS requirements. The IEEE 802.16 standard, however, leaves unstandardized the resource management and scheduling mechanisms, which are crucial components to guarantee QoS performance.

In this thesis, we have evaluated the performance of WiMAX networks in both fixed and highly mobile environments and tackled most of the resource management and scheduling issues that have been left open with the objective of defining an architecture that fulfills the QoS expectations of the five classes of applications addressed by the IEEE 802.16 standard.

In Chapter 1, we have provided an overview of the main features supported by PHY and MAC layers specified by the IEEE 802.16 standard, particularly focusing on the aspects useful for understanding our performance evaluation, carried out in Chapter 2. All the features related to QoS support at MAC level have been further discussed in Chapter 3. Indeed, because there are so many concepts to be introduced in this context, we have preferred to dedicate a whole chapter to this purpose.

In Chapter 2, an analytical framework was developed to investigate the performance bounds of OFDM-based 802.16 systems. This analytical framework was carried out with respect to what have been specified in the IEEE 802.16 standard [1]. It outlines a number of key features proposed by the standard and that have been hardly addressed in previous research works. Based on this framework, several scenarios were considered to evaluate the performance bounds of 802.16 systems under different MAC and PHY settings. The obtained results highlight the importance of considering the MAC and PHY overhead when evaluating the performance of IEEE 802.16 systems. Indeed this overhead, that is usually ignored or roughly estimated in most research works related to WiMAX resource allocation, may constitute 80 % of the whole frame. Also we have shown that using a larger bandwidth channel may yield minimal improvements on MAC performances. Also when investigating fragmentation and packing impact on MAC performance, we have shown that packing may considerably improve the resulting throughput especially for traffic carrying fixed-size packets.

As mentioned above, Chapter 3 was dedicated to introducing the features related to QoS support in WiMAX networks. More specifically, we aimed at providing a better understanding of the supported and missing features to ensure QoS support in WiMAX networks. Therefore, we have first described the main elements specified by the IEEE 802.16 standard to provide QoS for heterogeneous classes of traffic. Then, we have proposed a generic QoS framework that is independent of the adopted scheduling and CAC strategy. The proposed framework is intended to be a compilation of what we consider as key elements to handle QoS in WiMAX systems. We also addressed scheduling and admission control issues, highlighting the main challenges faced when designing a scheduling and/or CAC solution for WiMAX networks. These constraints represent also the main evaluation criteria for the different resource management mechanisms proposed in this work-in progress area.

---

The state of the art of these mechanisms was presented in Chapter 4 where we survey, classify, and compare different scheduling and CAC mechanisms proposed for WiMAX networks.

In Chapter 5, we have proposed a QoS architecture for PMP 802.16 systems operating in TDD mode over WirelessMAN-OFDM physical layer. It includes a CAC policy and a hierarchical scheduling algorithm. The proposed CAC policy adopts a Min-Max fairness approach making efficient and fair use of the available resources. The proposed scheduling algorithm flexibly adjusts uplink and downlink bandwidth to serve unbalanced traffic. This adaptive DL/UL scheduling procedure adapts the frame-by-frame allocations to the current needs of the connections with respect to the grants boundaries fixed by the CAC module. These boundaries may be set through a degradation of the ongoing connections rates if the available resources are not enough to accommodate the needs of a new connection for example. Through simulation, we reveal the efficiency of the proposed CAC scheme and show that our scheduling algorithm can meet the data rate requirements of the scheduling services specified by the IEEE 802.16 Standard. The degradation policy adopted in the proposed QoS solution can be handled by UDP traffic. However, it might cause an uneven behavior for TCP traffic especially under short round trip time (RTT) conditions. To prevent such a behavior, an extension of this work would be to combine our CAC policy with a TCP-friendly traffic policing mechanism among those available in the literature [56]. A further challenge we face would be to support bursty traffics and to integrate delay constraints in our proposal.

These two shortcomings were addressed in our multi-Constraints Scheduling Strategy (mCoSS) which is presented in chapter 6. mCoSS is designed for PMP 802.16 systems operating in TDD mode over OFDM or band-AMC OFDMA PHYs. Unlike the first QoS solution, mCoSS supposes the use of a predefined DL/UL ratio set by the operator. Most of the hierarchical scheduling strategies proposed in the literature and described in Chapter 4 (such as [13, 10, 9]) propose a specific queuing discipline for each scheduling service type, which increases significantly the complexity of the proposed scheduling policy. Unlike those approaches, the multi-Constraints Scheduling Strategy (mCoSS) proposed in Chapter 6 is designed to be applicable to all service types. Based on a modified dual-bucket traffic shaping mechanism used for all the scheduling service types, mCoSS allies the genericity of the approach to the specificity of the configuration since the dual-bucket mechanism is configured on a per-flow basis.

This shaping mechanism is combined with a two-rounds scheduling strategy which reflects (i) at the first round, the minimum data rates and latency requirements the BS or MS is committed to provide and (ii) at the second round, the efficiency and fairness of the resources management since the remaining bandwidth is shared in this round using a simple weighted fair queuing (WFQ) strategy; the allocations should nevertheless remain within the thresholds set by the dual-bucket shaping mechanism. The bandwidth request and grant mechanism adopted in mCoSS is designed to make a tradeoff between increasing the accuracy of the bandwidth needs perception at the BS and decreasing the overhead associated to frequent unicast polling. Indeed the proposed strategy alternates between bandwidth stealing, piggybacking, unicast, broadcast and group polling, and the use of PM bit according to the considered scheduling service type and the available resources. The proposed mCoSS has been implemented under QualNet 4.5 simulator and compared to strict-priority (SP) and to a variant of WFQ discipline. The preliminary results reported in this thesis validate and confirm the shaping fairness and AMC support capability of the proposed mCoSS. They also show that, compared to SP and WFQ, mCoSS provides better and more stable end-to-end-delay and jitter performances. The performance evaluation of mCoSS presented in Chapter 6 is by no means comprehensive. More simulation scenarios—involving more service types—need to be considered to check and validate other aspects of the proposed scheduling strategy.

We focus, in the last part of the thesis (Chapter 7 and Appendix A), on WiMAX technology from a mobility perspective. Several issues, such as horizontal and vertical handover support

---

in networks involving WiMAX systems, are studied and discussed in this part (Appendix A). A special emphasis has been put in Chapter 7 on evaluating the performance of Mobile WiMAX technology as a radio access technology (RAT) for intelligent transportation systems (ITS). Thus, we have investigated, through simulation, the potential and limitations of WiMAX as a communication media for vehicle-to-infrastructure (V2I) communications in comparison with the 5.9 GHz technology, based on the upcoming 802.11p standard. The performance of the two systems is evaluated for different vehicle speeds, traffic data rates, and network deployments. This comparative study is meant to be the first step towards defining optimal rules for choosing the most appropriate RAT among those proposed for next generation ITS.

---



## Appendix A

# Topics Related to Mobility Management in WiMAX Networks

The WiMAX forum estimates that more than 133 million of people will be using the WiMAX technology by the year 2012. From these users, more than 70% are expected to be using the mobile implementation of the technology. From this perspective, mobility management is a key aspect to provide access for these potential 70% of WiMAX users.

This appendix focuses on the latter topic. It describes the concepts and mechanisms introduced by the IEEE 802.16e standard—the amendment of the IEEE 802.16d-2004 standard—which provides enhancements mainly related to mobility management. We also cover, through this appendix, the main topics related to WiMAX networks from a mobility perspective and point out the research issues where there is room for contribution. The appendix is organized as follows. Section A.1 describes the logical architecture of a mobile WiMAX network. This architecture has been defined by the Network Working Group<sup>1</sup> (NWG) of the WiMAX Forum. Section A.2 describes the horizontal handoff procedure proposed by the IEEE 802.16e standard. Section A.3 presents some procedures, proposed in the literature, aiming at improving the handover mechanism.

Moreover, because this technology is more likely to co-exist with other access technologies in future networks, we dedicate Section A.4 to study the vertical handover mechanisms in heterogeneous environment involving mobile WiMAX systems. Roaming, which has been referred to as "the missing piece of the WiMAX puzzle", is briefly addressed in Section A.5. Section A.6 concludes the appendix by highlighting the main conclusions.

### A.1 Mobile WiMAX architecture

A Network Reference Model (NRM), presenting the logical architecture of a WiMAX network, has been proposed by the NWG [34]. It has been developed with the objective of supporting many architectural profiles and addressing multiple deployment scenarios of mobile WiMAX networks. In this section, we first describe the different entities of the NRM and then discuss the technical and business merits of each profile.

As shown in Figure A.1, the WiMAX NRM consists of three logical entities (Mobile Station MS, Access Service Network ASN, and Connectivity Service Network CSN) interconnected by R1-R5 reference points. These reference points insure multi-vendors interoperability between the different logical entities belonging to the network. Each of the MS, ASN, and CSN represents a

---

<sup>1</sup>A working group from the WiMAX forum. It is responsible for creating higher level networking specifications for fixed, nomad, portable and mobile WiMAX systems.

---

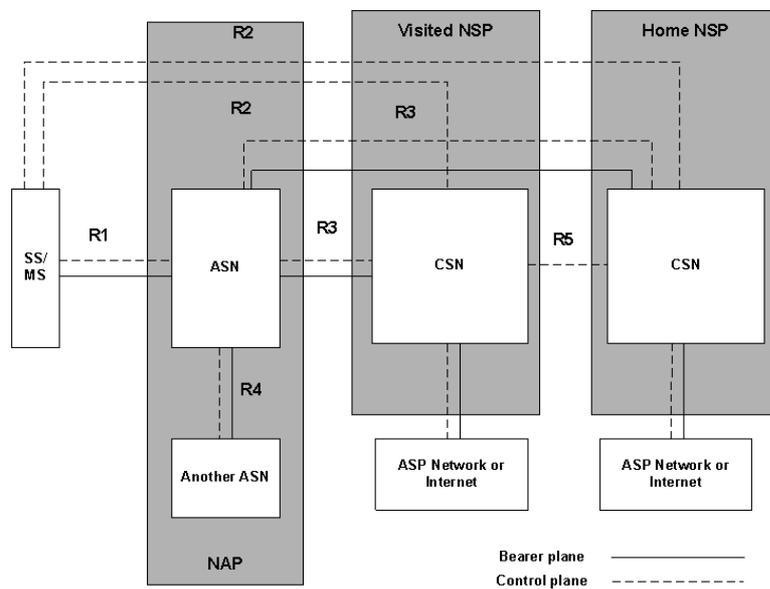


Figure A.1: Network Reference Model

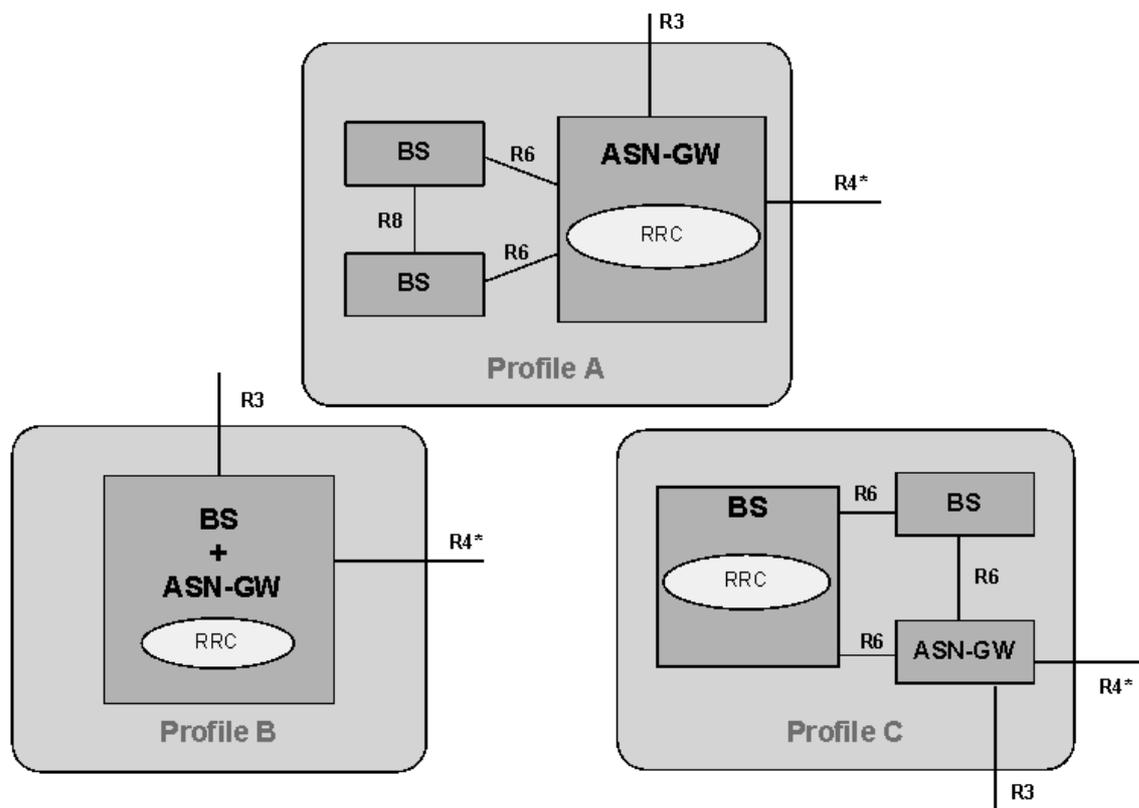


Figure A.2: ASN interoperability Profiles [33, 34]

grouping of functional entities (within an ASN, between an ASN and a MS, between an ASN and a CSN, etc.) that may be realized by a single or multiple physical devices:

1. Mobile Station (MS) is a generalized mobile equipment set which provides connectivity between a WiMAX subscriber equipment and a base station (BS).
2. Access Service Network (ASN) refers to a set of network functions providing radio access to the WiMAX MS. The mandatory functions that need to be provided by the ASN are: L2 and L3 connectivity with WiMAX subscriber, radio resource management (RRM), relay of AAA (Authentication, Authorization and Accounting) messages, network discovery and selection, mobility management, etc. An ASN consists of one or more BSs and one or more ASN-Gateways (ASN-GW):
  - (a) Base Station (BS) is a logical entity that incorporates a full instance of MAC and PHY layers compliant with the IEEE 802.16 suite of applicable standards.
  - (b) ASN-Gateway (ASN-GW) is a logical entity that represents an aggregation of control plane functions. It may also perform bearer plane routing or bridging function.
3. Connectivity service network (CSN) refers to a set of network functions that provide IP connectivity functions to the WiMAX subscribers. Among the functions that the CSN may provide, we find: Internet access, inter-ASN mobility, admission control based on user profiles, etc. The CSN may include network elements such as routers, AAA proxy/servers, user databases, etc.

The distribution of the different functions within the ASN (between the BS(s) and the ASN-GW(s)) is an implementation choice. Nevertheless, to guarantee network interoperability requirements, the NWG Release 1.0.0 [34] defines three different implementations of the ASN. These implementations, whose respective reference models are depicted in Figure A.2, are called interoperability profiles A, B, and C. Each of them corresponds to a specific distribution of ASN functions between the two entities composing the ASN: the ASN-GW(s) and the BS(s). As we can see it from Figure A.2, in Profile A, for instance, the radio resource control RRC (which is given here as example for function mapping) is in the ASN-GW while in Profile C it is accomplished by the ASN-GW. In Profile B, however, all the functions are located within a single ASN entity, which includes the case where all the functions are grouped in the same physical device. As discussed in [33], each profile has its own technical and business merits and selecting one or combining two or more of these profiles may seriously impact the handoff support in WiMAX networks. In [33], Hu *et al.* have investigated both the hierarchical and flat network architectures and their respective impacts on the performance of handoff in terms of latency, scalability, complexity, financial cost, etc. The authors have then mapped the different interoperability profiles to a hierarchical, flat or hybrid design which could help to choose the most appropriate architecture when deploying a technological solution.

## A.2 Horizontal handover in 802.16e

The IEEE 802.16e standard [2] defines three handover schemes:

- a mandatory hard HO mode also known as break-before-make HO. In this mode, the air interface link between the MS and the Serving BS is broken at all layers before being established again at the target BS. The HO process may be initiated either by the MS or by the BS.

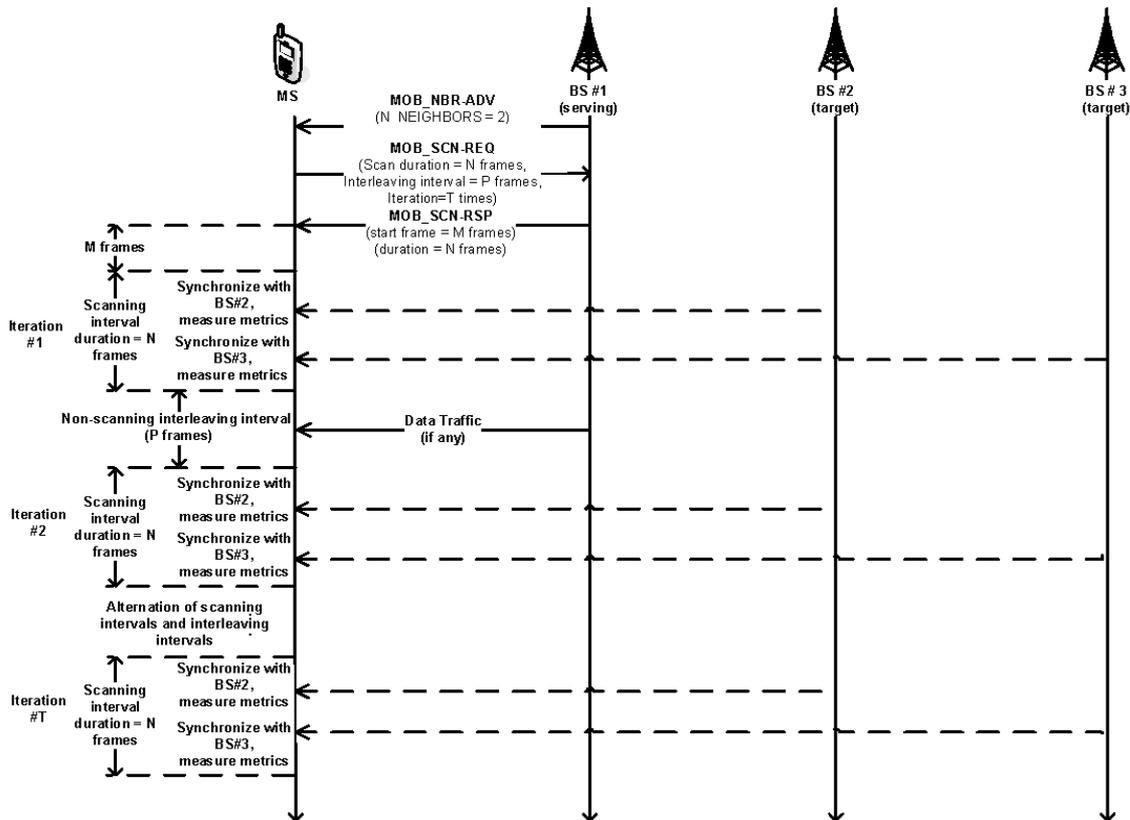


Figure A.3: Example of neighbor BS advertisement and scanning (without association) by MS request [2]

- two optional soft, known as make-before-break, HO modes:
  - Macro diversity HO (MDHO): this mode is defined in [2] as the process in which the MS migrates from an air interface provided by one or more BSs to the air interface provided by one or more other BSs. In the DL (respectively UL), this is achieved by having two or more BSs transmitting (respectively receiving) the same PDU to (respectively from) the MS.
  - Fast BS switching (FBSS): in this mode, an active set is maintained. It consists of a set of candidate BSs to which the MS is likely to handoff in near future. At any given frame, the MS is exchanging data only with one BS—anchor BS—of this active set [2].

More details about these three modes are provided in this section. Nevertheless more insight is given on the hard HO mode which is the only mandatory mode.

### A.2.1 Network topology acquisition

1. The BS periodically broadcasts the network topology information using the MOB\_NBR-ADV message. The message includes the BSIDs of the neighboring BSs along with their respective channel characteristics normally provided by each BS own Downlink/Uplink Channel Descriptor (DCD/UCD) message transmission. This information is intended to

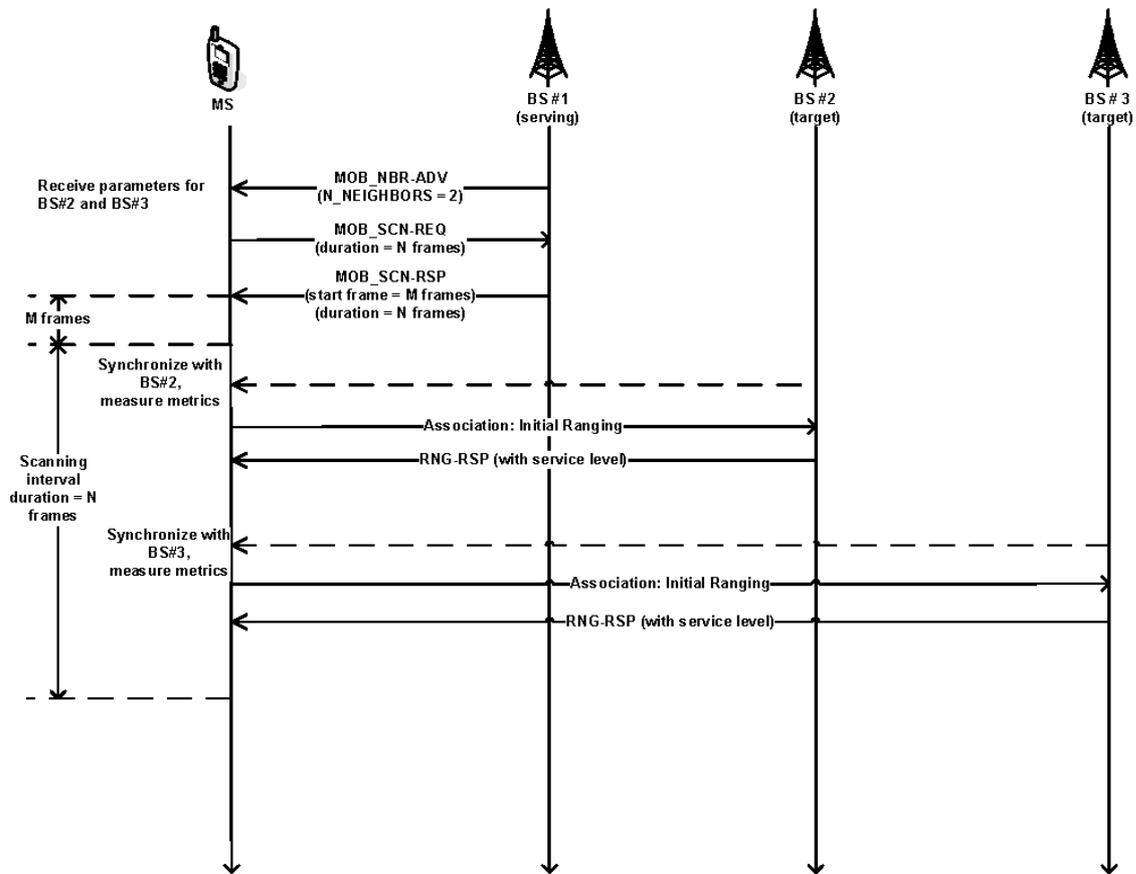


Figure A.4: Example of neighbor BS advertisement and scanning (with non-coordinated association) by MS request [2]

enable the MS to perform fast synchronization with the advertised BSs by removing the need to monitor the DCD/UCD broadcasts from each neighboring BS.

2. Based on the information provided by the MOB\_NBR-ADV, the MS becomes aware of the neighboring BSs and triggers the scanning and synchronization phase. Indeed, to handoff, the MS needs to seek available BSs and check if they are suitable as possible target BSs. Therefore, the MS sends MOB\_SCN-REQ message to the serving BS indicating a group of neighboring BSs for which a group of scanning intervals is requested. The MOB\_SCN-REQ message includes the requested scanning interval duration, the duration of the interleaving interval, and the requested number of scanning iterations. In the example illustrated in Figure A.3, these parameters correspond to P frames, N frames, and T iterations, respectively. Note that the scanning phase could be triggered by the serving BS. If it is the case, the serving BS shall send to the MS a MOB\_SCN-RSP message indicating a list of recommended neighboring BSs.
3. Upon reception of the MOB\_SCN-REQ message, the serving BS shall respond with a MOB\_SCN-RSP message. In this message, the serving BS either grants a scanning interval at least as long as the one requested by the MS (which is the case in our example A.3) or rejects the request.
4. After receiving the MOB\_SCN-RSP message granting the request, the MS may scan—beginning at *Start frame*—one or more BSs during the time allocated by the serving BS. Each time a neighboring BS is detected through scanning, the MS may attempt to synchronize with its downlink transmissions and estimate the quality of the PHY channel to evaluate its suitability as a potential target BS in the future. The serving BS may ask (by setting the report mode field to 0b10 in the MOB\_SCN-RSP) the MS to report the scanning results by transmitting a MOB\_SCN-REP.
5. During the scanning interval, the serving BS may buffer incoming data addressed to the MS and then transmit that data during any interleaving interval after the MS has exited the scanning mode.

Depending on the value of the scanning type field indicated in the MOB\_SCN-REQ, the MS may request either scanning only or scanning with association. The association procedure is an optional ranging phase that may be performed during the scanning interval. It enables the MS to acquire and record ranging parameters—by adjusting the time offset, the frequency and the power level—to be used to choose a potential target BS. The standard IEEE 802.16e [2] defines three levels of association:

- Association Level 0 — scan/association without coordination: the target BS has no knowledge of the scanning MS and only provides contention-based ranging allocations.
  - Association Level 1 — association with coordination: the serving BS coordinates the association between the MS and the requested neighboring BSs. Each neighbor (NBR) BS provides a ranging region for association at a predefined “rendezvous time” (corresponding to a relative frame number). It also reserves a unique initial ranging code and a ranging slot within the allocated region. The NBR BS may assign the same code or ranging slot to more than one BS but not both, so that no potential collision may occur between transmissions of different MSs.
-

- Association Level 2 — network-assisted association reporting: the procedure is similar to level 1 except that the MS does not need to wait for RNG-RSP from the NBR BS. The ranging response is sent by the NBR BS to the serving BS over the backbone, which then forwards it to the MS.

## **A.2.2 Handover process**

The handover is defined as the process in which a MS migrates from the air-interface provided by one BS (the serving BS) to the air-interface of another BS (target BS) [2]. It consists of the following phases:

### **A.2.2.1 Cell reselection**

Cell reselection refers to the process of an MS Scanning and/or Association with one or more BSs (as described in Section A.2.1) in order to determine their suitability, along with other performance considerations as a handover target [2]. The information acquired from the MOB\_NBR-ADV message might be used by the MS to give insight into available neighboring BSs for cell reselection considerations.

### **A.2.2.2 HO decision and initiation**

The handover process begins with a decision that originates either from the MS, or the BS (the BS can force the MS to conduct handover), or on the network. A handover could be decided for many reasons; for example when the MS performance at a potential target BS is expected to be higher than at the serving BS. Note that the handover decision algorithm is beyond the scope of 802.16e standard, which leaves room for research contributions.

Once a handover is decided, it is notified through a MOB\_MSHO-REQ or a MOB\_BSHO-REQ indicating one or more possible target BSs. If the handover request is formulated by the MS, it shall be acknowledged with a MOB\_BSHO-RSP. When the handover is initiated by the BS, it could be either recommended or mandatory. If it is a mandatory handover, the MS shall send MOB\_HO-IND to the serving BS. The MOB\_HO-IND may indicate a HO reject when the MS is unable to handoff to any of the recommended target BSs listed in the MOB\_BSHO-REQ.

### **A.2.2.3 Synchronization to target BS downlink**

MS shall synchronize to downlink transmissions of target BS and obtain DL and UL transmission parameters. This process may be shortened in two cases: (i) if the MS had previously received a MOB\_NBR-ADV message including target BSID, physical frequency, DCD and UCD, or (ii) if the target BS had previously received HO notification from serving BS over the backbone in which case the target BS may allocate a non-contention-based initial ranging opportunity for the MS.

### **A.2.2.4 Ranging and network re-entry**

After adjusting all the PHY parameters, the network re-entry process is initiated between the MS and the target BS. The network re-entry procedure normally includes the following steps (i-iv).

(i) Negotiation of basic capabilities: the MS and the target BS exchange their supported parameters such as the current transmit power or the security parameters support. This step is performed by exchanging SBC-REQ and SBC-RSP management messages.

---

(ii) Privacy key management (PKM) authentication phase: during this phase, the MS exchanges secure keys with the target BS. The MS sends a PKM-REQ message and the BS responds with a PKM-RSP message.

(iii) Traffic encryption keys (TEK) establishment phase.

(iv) Registration: the registration is the process by which the SS is allowed to enter into the network [2]. The registration is performed by exchanging REG-REQ and REG-RSP between the MS and the target BS.

The network re-entry process may be shortened since the target BS may decide to skip one or more of these steps (i-iv) if it disposes of the corresponding information obtained from the serving BS over the backbone.

#### **A.2.2.5 Termination of MS context**

The termination of the MS context is the final stage of the handover procedure. In this step, the serving BS proceeds to the termination of all the connections belonging to the MS along with their associated context (information in the queues, timers, counters, etc.).

Note that the handover procedure might be canceled by the MS at any time prior to the expiration of Resource\_Retain\_Time interval after transmission of MOB\_HO-IND message.

### **A.2.3 Fast BS switching (FBSS) and macro diversity handover (MDHO)**

As mentioned before, in addition to the hard handover procedure previously described, the IEEE 802.16e standard defines two optional handover modes: MDHO and FBSS. The MDHO or FBSS capability can be enabled or disabled in the REG-REQ/RSP message exchange. In both modes, a Diversity Set is maintained. The Diversity Set is a list of selected BSs that are involved in the MDHO or FBSS process. These BSs should be synchronized in both time and frequency and are required to share the MAC context associated to the MS. The MAC context includes the parameters that are normally exchanged during the network entry along with the service flows associated to the MS connections.

#### **A.2.3.1 Macro diversity handover (MDHO)**

A MDHO begins with a decision for an MS to transmit to and receive from multiple BSs at the same time. This decision is communicated through MOB\_BSHO-REQ or MOB\_MSHO-REQ messages. When operating in MDHO mode, the MS communicates with all the BSs belonging to the Diversity Set for DL and UL unicast messages and traffic. For DL MDHO, two or more BSs provide synchronized transmission of MS data so that the MS performs diversity combining. For UL MDHO, the MS data transmission is received by multiple BSs so that they can perform selection diversity of the received information.

To monitor DL control information and DL broadcast messages, the MS can use one of the following two methods. The first method is the MS monitors only the Anchor BS—a BS defined among the Diversity Set—for DL control information and DL broadcast messages. In this case, the DL-MAP and UL-MAP of the Anchor BS may contain burst allocation information for the non-Anchor Active BS. The second method is the MS monitors all the BSs in the Diversity Set for DL control information and DL broadcast messages. In this case, the DL-MAP and UL-MAP of any Active BS may contain burst allocation information for the other Active BSs. The method to be used by the MS is defined during the REG-REQ and REG-RSP handshake.

---

### A.2.3.2 Fast BS switching (FBSS)

FBSS HO begins with a decision for an MS to receive/transmit data from/to the Anchor BS that may change within the Diversity Set. A FBSS can start with MOB\_BSHO-REQ or MOB\_MSHO-REQ messages. When operating in FBSS mode, the MS is required to continuously monitor the signal strength of the BSs belonging to the Diversity Set. The MS shall select a BS from its current Diversity Set to be the Anchor BS and report the selected Anchor BS on MOB\_MSHO-REQ message. BS switching i.e. transition from the Anchor BS to another BS is performed without invocation of the handover procedure described in Section A.2.2.

The BS supporting MDHO or FBSS shall broadcast the DCD message that includes the H\_Add Threshold and H\_Delete Threshold. These thresholds are used by the FBSS/MDHO capable MS to determine if MOB\_MSHO-REQ should be sent. When long-term CINR of a BS is less than H\_Delete Threshold, the MS shall send MOB\_MSHO-REQ to require dropping this BS from the Diversity Set. When long-term CINR of a neighbor BS is higher than H\_Add Threshold, the MS shall send MOB\_MSHO-REQ to require adding this neighbor BS to the diversity set. Figure A.5 illustrates an example of a Diversity Set update—add of a new BS—during a MDHO procedure.

## DISCUSSION

From the description of the three handover modes, the hard handoff procedure consists of more steps and might cause intolerable delays for real-time traffic. Nevertheless, the two soft handover modes FBSS and MDHO cannot be a reliable alternative to the mandatory hard HO scheme for many reasons. On the one hand, as we have mentioned before, there are several restrictions on BSs working in MDHO/FBSS modes since they need to synchronize on time (same time source) and frequency and have synchronized frame structures which entails extra costs. On the other hand, in both FBSS and MDHO modes, the BSs in the same Diversity Set are likely to belong to the same subnet while a handover may occur between BSs in different subnets. Therefore, in the remaining of the appendix, more insight will be given into the hard handover scheme. More specifically, we will present some works aiming at optimizing the hard handover procedure in IEEE 802.16e networks.

## A.3 Optimized 802.16e handover schemes

Improving the handoff process in mobile WiMAX networks is a topic that have received a lot of attention in the last few years. Indeed, in order to enable always-on connectivity, it is necessary to achieve a fast and smooth handoff over the network. To reach that goal, the research works addressing this issue have adopted mainly two approaches: improving the handover at Layer 2 or considering a cross-layer mechanism in which L2 and L3 collaborate to have better results.

### A.3.1 L2 handover schemes

In order to reduce the handover delay, Lee *et al.* [63] have focused on eliminating the redundant processes existing in the handover procedure defined in the IEEE 802.16e standard [2]. The approach consists in using a target BS estimation algorithm to select a HO target BS instead of scanning, one by one all the neighboring BSs. The target BS estimation algorithm assumes that the NBR BS with bigger mean CINR and smaller arrival time difference is more likely to be the target BS. The MS does not need then to associate to the neighboring BSs. However, we consider that, by eliminating both the scanning and the association phases, the handover decision loses its accuracy since the MS does not dispose of information precise enough to make a handoff decision.

---

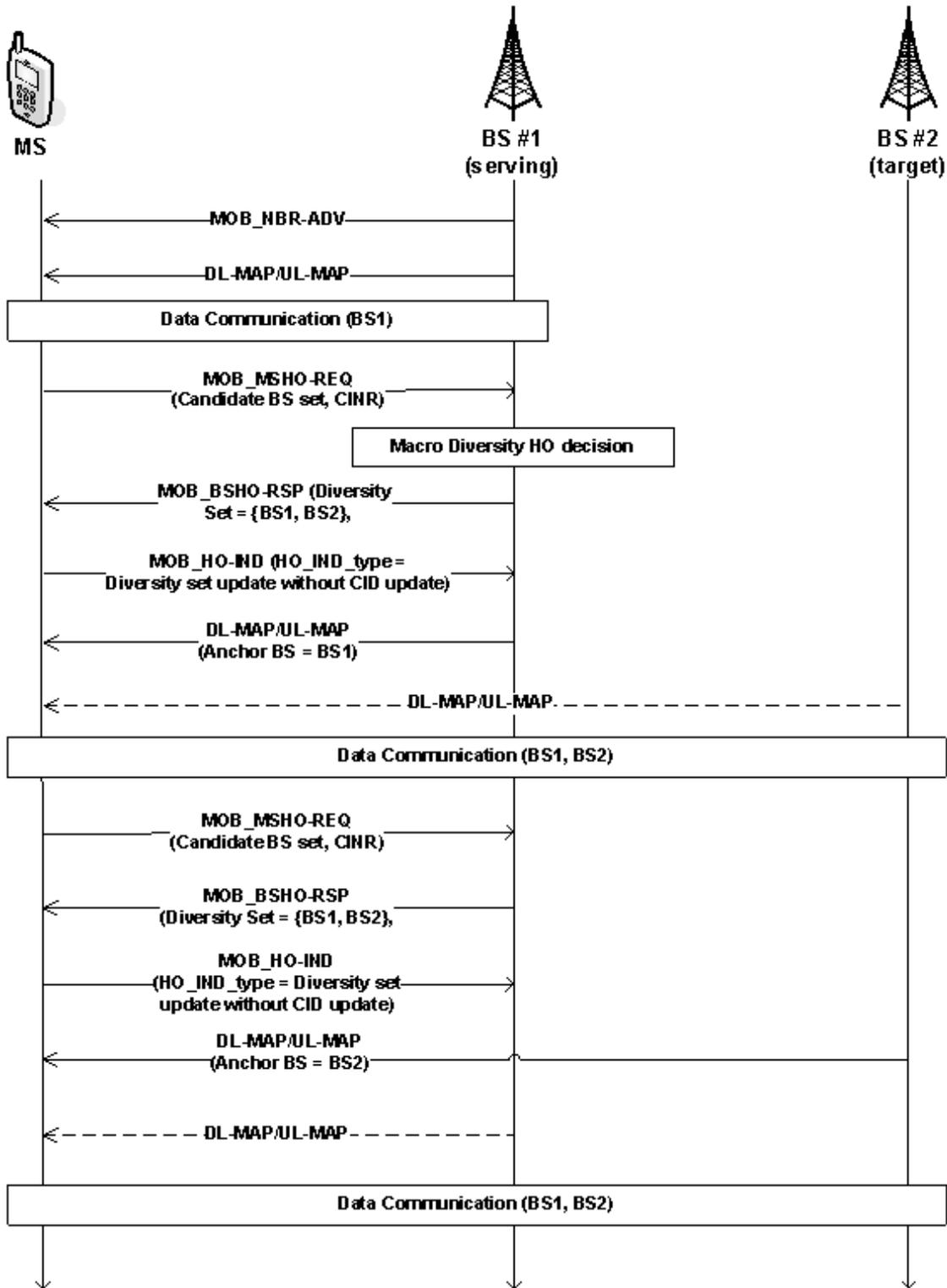


Figure A.5: Example of macro diversity HO (Diversity Set Update: Add) [2]

---

Instead of predicting the potential target BS, Chen *et al.* propose in [64] a pre-coordinated handover mechanism in which the handover time is predicted. In the proposed mechanism, the distance between the MS and the serving BS is calculated to estimate the needed handover time, then a pre-coordination is performed with the target BS.

In order to locate the position of the MS, the serving BS measures the signal-to-noise ratio (SNR) of the mobile station every 10 s. Based on that, the distance, the direction, and the velocity of the MS are derived. If the MS approaches the boundary  $h$  of the serving base station macrocell, the serving BS pre-coordinates a handover with the "only" target BS in that direction. The pre-coordination phase consists in sending a MOB\_BSHO-REQ to the target BS which would respond with a MOB\_BSHO-RSP in which it allocates—if it has enough resources—a fast ranging opportunity for the MS and specifies its PHY parameters. The target BS will have then to hold this request service for 10 s. When the MS requests to handoff (estimated to 10 s before the predicted handover time), the BS responds by MOB\_NBR-ADV message in which it includes the information transmitted by the target BS. This would facilitate the migration of the MS to the new channel and thus reduce the disruption time. Nevertheless, the performance estimation algorithm needs further investigation to be reliable.

### A.3.2 L2-L3 cross-layer handover schemes

In [65], Chen *et al.* have proposed a cross-layer handover scheme in which they use layer 3 to transmit MAC messages between the MS and the BSs (the serving BS and the NBR BSs) during the handoff process. In the proposed cross-layer scheme, two tunnels are created to redirect and relay these messages: an L2 tunnel between the MS and the serving BS and an L3 tunnel between the serving and the neighboring (target) BSs. The idea behind the creation of these tunnels is to minimize the delay due to direct messages transportation between the MS and NBR BSs which constitutes a source of latency in the handover process. When the handover is requested, the serving BS negotiates for the MS a fast ranging opportunity from the neighboring BSs. The MS then switches to the channels to be scanned and tries to synchronize with each associated NBR BSs. Once the synchronization is performed, the MS sends a MOB\_RNG-REQ on each channel. However, unlike the regular handover procedure described in Section A.2, the MS does not need to wait for RNG-RSP from each scanned NBR BS. Instead, the MS informs the BS that the ranging request phase has finished by sending a RNG\_RSP-REQ message (a new management message proposed by Chen *et al.* [65]) and restores the uplink transmission. Upon reception of the RNG\_RSP-REQ, the serving BS understands that the MS is ready to receive the RNG\_RSP messages. These messages have been encapsulated by the NBR BSs and sent to the serving BS which decapsulates and stores them before forwarding them to the MS. This way, the uplink transmission is restored faster.

Moreover, a fast re-entry procedure is proposed. Instead of disconnecting and connecting with the target BS as described in Section A.2, the MS sends all the messages to the serving BS which relays them to the target BS through the IP backbone.

The idea of combining L2 and L3 mechanisms to shorten the handover time and to allow handover between different subnets has been also investigated by Chang *et al.* in [66]. The authors have mainly focused on interleaving the authentication process with a fast handover mechanism to speedup the handover process while securing the whole mechanism. Chang *et al.* have based their proposal on a draft version of an RFC [35]—recently finalized by IETF—proposing Mobile IP fast handover mechanism over IEEE 802.16e networks.

---

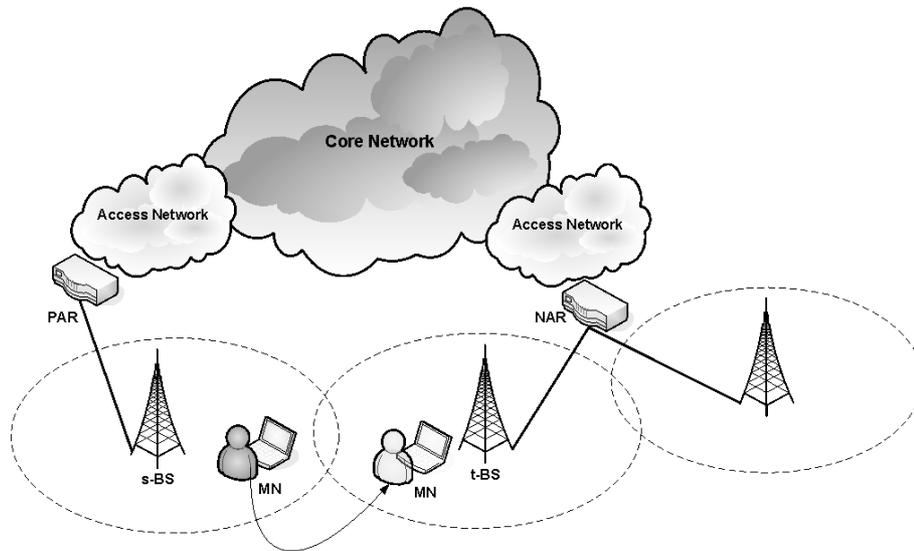


Figure A.6: Example of a handover between two different subnets

### A.3.3 Mobile IPv6 fast handovers over IEEE 802.16e networks

This section is dedicated to the description of the interleaving between 802.16e and fast mobile IPv6 (FMIPv6) handover mechanisms proposed by IETF in [35]. The handoff procedure is explained through two examples corresponding to the predictive (Figure A.7) and reactive mode (Figure A.8), respectively.

#### A.3.3.1 Predictive mode

The different steps commented in this section are illustrated in Figure A.7.

##### Access router discovery

- 1-3 When a new BS (Point of Attachment PoA) is detected through the reception of MOB\_NBR-ADV or through scanning, the link layer of the MS triggers a NEW\_LINK\_DETECTED primitive to the IP layer.
4. When receiving the NEW\_LINK\_DETECTED from the link layer, the IP layer sends a router solicitation message RtSolPr (Router Solicitation for Proxy Advertisement) to the previous access router (PAR) to acquire the L3 parameters of the access router associated to the new PoA (the new BS). The PAR responds by sending a Proxy Router Advertisement (PrRtAdv) that provides information such as the router address and additional parameters about neighboring links.

The objective of this step is to enable the quick discovery—in IP layer—of the access router associated to the new BS.

##### Handover preparation

5. When the MN decides to change the PoA (because of a degradation in signal strength, or for better QoS, etc.) it initiates a handover procedure by sending a MOB\_MSHO-REQ to the serving BS which will respond by a MOB\_MSHO-RSP. As we have seen in Section A.2.2, the handover might also be initiated by the serving base station (MOB\_BSHO-REQ).

6. Once a MOB\_MSHO-RSP/MOB\_BSHO-REQ is received, the link layer triggers a LINK\_HAN-DOVER\_IMPENDING primitive, enclosing the decided target BS, to inform the IP layer that a link layer handover decision has been made and that its execution is imminent.  
Based on the information collected during the access router discovery phase, the IP layer checks whether the target BS belongs to a different subnet (c.f. Figure A.6). If the target network proves to be in the same subnet, the MN can continue to use the same IP address and thus, there is no need to perform FMIPv6. Otherwise,
7. based on the information provided by the PrRtAdv, the IP layer formulates a prospective NCoA (New Care of Address) and sends a Fast Binding Update (FBU) message to the PAR. When received successfully, the FBU is processed by the PAR and the NAR according to RFC 5268 (FMIPv6 [67]).  
The PAR sets up a tunnel between the PCoA (Previous Care of Address) and the NCoA by exchanging a HI (Handover Initiation) and HAcK (Handover Acknowledgment) messages with the NAR. In the HAcK message, the NCoA is either confirmed or re-assigned by the NAR. Finally, the NCoA is transmitted to the MN through the FBack (Fast Binding Acknowledgment) message in case of predictive mode (shown in Figure A.7) and the packets destined to the MN are forwarded to the NCoA. The difference with the reactive mode will be explained at the end of this section.

#### **Handover execution**

8. If the MN receives a FBack on the previous link, it sends a MOB\_HO-IND message as a final indication of handover. Optionally, the LINK\_SWITCH command could be issued by the IP layer upon the reception of FBack to force the MN to switch from an old BS to a new BS. This command forces the use of predictive mode even after switching to the new link.
9. Once the links are switched, the MN synchronizes with the new PoA (target BS) and performs the 802.16e network entry procedure. As we have mentioned before in Section A.2.2, this phase (or some of its steps) might be omitted if the serving BS had transferred the MN context to the target BS over the backbone.
10. Once the network entry is completed, the link layer triggers a LINK\_UP primitive to inform the IP layer that it is ready for data transmission.

#### **Handover completion**

10. When the MN IP layer receives the LINK\_UP primitive, it checks whether the target network is the one predicted by the FMIPv6 operation. If it is the case, it sends an Unsolicited Neighbor Advertisement (UNA) message to the NAR (predictive mode) using the NCoA as source IP address and starts performing the DAD (Duplicate Address Detection) for the NCoA.
11. As soon as the UNA message is received, the NAR transfers the buffered packets to the MN.

#### **A.3.3.2 Reactive mode**

The different steps commented in this section are illustrated in Figure A.8.

---

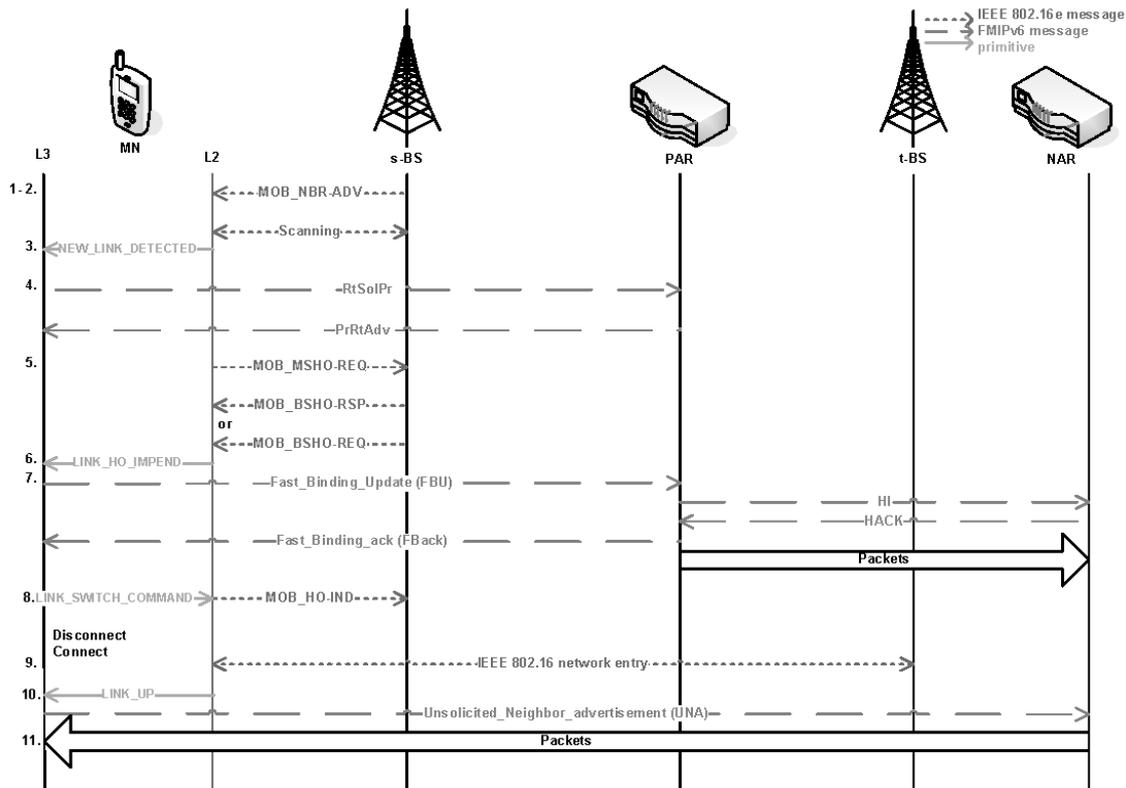


Figure A.7: Predictive fast handover in 802.16e [35]

### Access router discovery

1~4 The same procedure as in predictive mode

### Handover preparation

5~7. The same procedure as in predictive mode. Nevertheless, note that the FBU has not reached the PAR, and so no FBack has been received by the MN either.

8. Unlike in predictive mode, the MN issues a MOB\_HO-IND without waiting for an FBack message. When receiving this final indication of handover (MOB\_HO-IND), the serving BS releases all the MN context which means that data packet transfer is no longer allowed between the MN and the BS (as we can see from Figure A.8).

### Handover execution

9. The MN conducts handover to the target BS and performs the 802.16e network entry procedure.

10. The MN link layer triggers a LINK\_UP primitive to inform the IP layer that it is ready for data transmission.

### Handover Completion

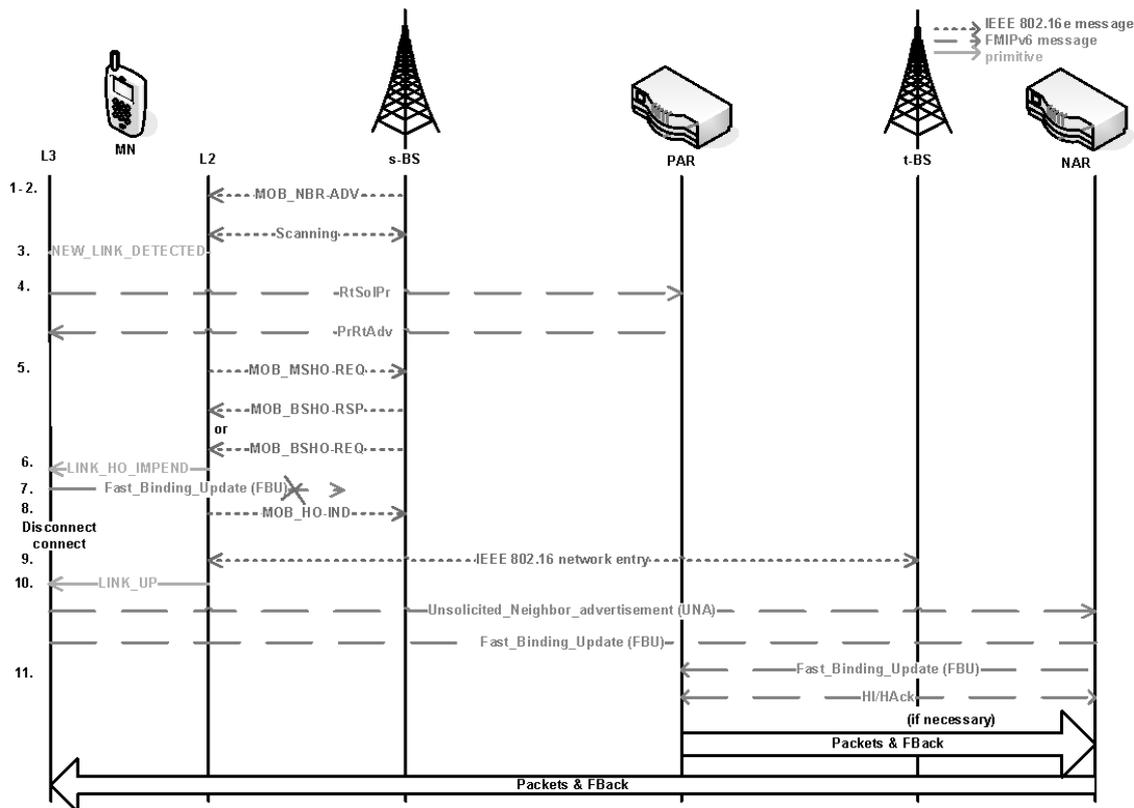


Figure A.8: Reactive fast handover in 802.16e [35]

10. Recall that, in reactive mode, the MN has moved to the target network without receiving an FBack message in the previous link. Therefore, upon reception of the LINK\_UP primitive, the IP layer sends (i) an UNA message to the NAR using the NCoA as source IP address to announce a link layer address change, and (ii) a FBU message to instruct the PAR to redirect its traffic towards the NAR.
11. When the NAR receives the UNA and the FBU from the MN, it exchanges a HI/HACK with the PAR. The FBack and Packets are then forwarded from the PAR and delivered to the MN through the NAR using the NCoA as destination IP address.

## Discussion

Mobile IPv6 fast handovers, like all cross-layer handover management mechanisms in general, are based on the collaboration of different layers in order to enhance the mobility management. This idea of integrating information from different network layers helps to improve the HO management performance. Nevertheless, because these solutions usually require significant modifications in the network stack, their deployment becomes prohibitive [68].

## A.4 Vertical handover

Next generation networks will more likely consist of heterogeneous networks such as integrated WiFi/WiMAX networks, WiMAX/CDMA2000 or networks combining WiMAX and 3G/4G tech-

nology. In this section, we describe the deployment of such hybrid networks and discuss the main challenging issues that arise when inter-networking WiMAX and other technologies. We focus on the vertical handover mechanisms proposed in the literature to guarantee the service continuity without QoS degradation for users switching from one network to another. The second part of this section is dedicated to the media-independent handover (MIH) mechanism proposed by the IEEE 802.21 task group. The recently published IEEE 802.21-2008 Standard [36] enables handover and interoperability between heterogeneous network types including both 802 and cellular networks.

#### **A.4.1 Vertical handover mechanisms involving 802.16e networks**

For both horizontal and vertical handover, the main objective is to provide a fast and seamless handover. However, because of the heterogeneity of the networks involved in the vertical handoff process, ensuring a continuous connectivity is even more challenging.

To make a horizontal handover decision, considering only the radio signal strength was enough while in a hybrid network environment this metric is not sufficient. Indeed, more parameters need to be considered: available bandwidth, latency, packet error rate, monetary cost, power consumption, user preferences, etc. [69].

In this section, we present works that have investigated the vertical handoff mechanisms involving mobile WiMAX networks. Each of these works have focused on the enhancement of one or more of the three main phases of a vertical handover procedure which are:

1. Finding candidate networks: also referred to as system discovery phase during which the MS needs to know which networks can be used.
2. Deciding a handoff: during this phase, the MS needs to evaluate the reachable wireless networks and to decide whether to keep using the same network or to switch to another network. This decision could involve several criteria: the type of applications running, their QoS requirements, the access cost, etc. [70].
3. Executing a handoff: a critical phase during which the connections need to be rerouted in a seamless manner with transfer of the user's context.

According to another classification proposed in [36], the two first steps could be merged into a single phase called "handover initiation" which encloses network discovery, network selection, and handover negotiation. Based on the same classification [36], the handover execution would correspond to two steps: handover preparation (L2 and L3 connectivity) and handover execution (connection transfer).

Whatever is the adopted classification, we notice that the phase on which most of the works have focused is the handover decision phase. In [71] for example, Dai *et al.* have proposed the use of two triggers: (i) connectivity trigger and (ii) performance trigger based on which the handoff between WiFi and WiMAX is decided. The first trigger is based on SINR indication to evaluate the risk of connection loss and would decide a handover if the SINR is below a certain SINR target and if other networks are detected. The performance trigger however, combines data rate and channel occupancy to derive an estimation of the current throughput and decide a potential handoff when needed (i.e. when the throughput is below a certain threshold).

In [72], the handover decision might be initiated either (i) by the user when it is moving and needs to gain in performance or (ii) by the WiMAX network to release resources and accommodate new calls (WiMAX calls) or VHO calls (from an UMTS network). The vertical handoff decision algorithm (VHDA) proposed in [72] depends of the improvement that could be gained from the

---



**MIH Function (MIHF)** The main role of the MIHF is to assist the network selector entity in making an effective network selection by providing all the necessary inputs for such a decision: QoS requirements, battery life constraints, monetary cost, user preferences, operators' policies, etc.. These information are meant to facilitate the handover decision and to maximize its efficiency. To achieve this role, the MIHF communicates with lower layers through technology-specific interfaces and provides services to the upper layers (MIH users) in a unified and abstracted way. More details about the services provided by the MIHF are given in Section A.4.2.2.

**MIH User (MIHU)** MIH users (MIHUs) are the entities responsible for mobility management and handover decision making. They reside at Layer 3 or above in the network stack. As examples of MIH users, we can cite MIP at network layer, mobile Stream Control Transmission Protocol (mSCTP) at transport layer, and Session Initiation Protocol (SIP) at application layer [73]. The MIHU base their handover decisions on their own internal policy but also on the information provided by the MIHF.

**SAPs** In order to make possible the communication between the different architectural components of the MIH framework, the IEEE 802.21 standard defines a set of SAPs with their associated primitives. Figure A.9 shows the different SAPs interfacing the MIHF with other layers:

1. The media-independent SAP MIH\_SAP allows the MIH users to access the MIHF services.
2. The link-layer SAPs MIH\_LINK\_SAP are media-dependent SAPs that allow the MIHF to gather link information and control link behavior during handovers. Each link-layer technology (e.g 802.3, 802.16, 3GPP, etc.) specifies its own technology-dependent SAPs and the MIH\_LINK\_SAP maps to these technology-specific SAPs. As example of media-specific SAPs, we can cite the C\_SAP, M\_SAP, and CS\_SAP that are defined in IEEE Std 802.16 to provide interfaces between the MIHF and different components of the 802.16 network stack; namely with the control plane (C\_SAP), the management plane functions (M\_SAP), and the service-specific Convergence Sublayer (CS\_SAP).
3. The MIH\_NET\_SAP is another media-dependent SAP that provides transport services over the data plane and allows the MIHF to communicate with remote MIHFs.

#### A.4.2.2 MIHF services

In order to facilitate the handover procedure across heterogeneous networks, the MIHF entity provides the three following categories of services to the MIH users: MIH information service (MIIS), MIH event service (MIES), and MIH command service (MICS).

**MIIS: MIH information service** The media independent information service allows the MIH users to acquire a general view about the networks present in the vicinity of the MN in order to enable a more effective handover decision. These information include for instance the list of available networks, their link-layer static information (e.g. whether QoS and security are supported in a particular network), and other geographical positioning information that could be used further to optimize the handover decision.

---

**MIES: MIH event service** Unlike the MIIS which provides a static (or rarely changing) information about the surrounding networks, the MIH event service (MIES) triggers dynamic changes in link conditions. Indeed, it provides event reporting about MAC and PHY state changes through triggers that indicate for instance that the L2 connection is broken (LINK\_DOWN) or that the link conditions are degrading and the loss of connectivity is imminent (LINK\_GOING\_DOWN). Other triggers might report the failure/success of PDUs transmission (e.g. Link\_PDU\_Transmit\_Status), or the handover status (e.g. Link\_Handover\_Complete).

**MICS: MIH command service** The MIH command service (MICS) refers to the set of commands that originate (i) either from the MIH users: MIH commands, (ii) or from the MIHF: link commands and are directed to the lower layers. MIH\_MN\_HO\_Candidate\_Query is an example of a remote MIH command used by the MN to query and obtain handover related information about possible candidate networks. The link commands are local commands that are used to control and configure the link layers (e.g. Link\_Configure\_Thresholds which is used to set link parameter thresholds) or to retrieve link-specific information (e.g. Link\_Get\_Parameters commands provide information about the SNR, the bit-error-rate BER, etc.)

**Service management** In order to benefit from the services provided by the MIHF, the MIH entities need to be configured properly using the following service management functions:

- MIH capability discovery

This step is necessary to the MN to discover local and/or remote MIHF capabilities in terms of MIH supported services. This could be performed either through the MIH protocol or via media-specific mechanisms (e.g. beacon frames for 802.11). For 802.16 networks for instance, the MN can use the management messages such as downlink channel descriptor (DCD), or uplink channel descriptor (UCD) to retrieve such information.

- MIH registration

MIH registration is defined to query access to certain MIH services. This phase is either mandatory or optional depending on the required level of service support. Indeed, the registration allows the peer MIHF entities to communicate in a trusted manner and gives them access to extensive information [74]. Nevertheless, for security issues, this registration is valid only for a certain period of time and has to be re-established when needed.

- MIH event subscription

refers to the fact of subscribing to a particular set of events that are provided by the MIES of a local or remote MIHF. By subscribing to a set of events and commands, the MIHU expresses for example its interest in triggering specific link behavior. Each subscription request needs to be individually validated by a confirmation from the event source (e.g. the peer MIHF) [74].

## Discussion

Because next generation networks will more likely consist of heterogeneous networks, the convergence towards a unified handover mechanism has become a must. From that perspective, the MIH mechanism offers an interesting alternative since it provides a generalized and standardized solution for handover across different access technologies. Nevertheless, its success highly relies on vendors support and willingness to integrate it in their future products [74].

---

## A.5 Roaming

Roaming is the process through which a mobile user automatically gains access to the services of a different provider, when outside the coverage area of its home network provider. Roaming service is made possible through Network Service Providers (NSPs) that have cooperative agreements to grant each others' customers local access to their resources. The WiMAX roaming relationship between NSPs consists in a technical and a business relationship.

Roaming provides significant advantages to customers, Home Network Service Providers (HNSP) and Visited Network Service Provider (VNSP) network operators. First, for users, they are able to use the network services even when traveling outside the coverage area of their HNSP. All the connectivity problems are transparent to them. From the HNSP point of view, roaming represents an increasing in the coverage footprint without incurring additional network capital costs. For the VNSP, roaming may provide additional revenue opportunities.

The roaming process may be considered outbound or inbound. For the HNSP a roaming is an outbound roaming, since the node is using the services of another operator. For the VNSP, it is an inbound roaming, since it is a user from another operator that is requesting to use the VNSP network.

Roaming can also be classified into national and international. National roaming occurs when the visited network is in the same country as the home network. International roaming occurs when the visited network is in a different country than the home network. Roaming can also occur between networks using different technologies, inter-standard roaming (which is referred to in this appendix as vertical handover), e.g. WiMAX and WiFi or WiMAX and GSM/CDMA [75].

To allow a more generic and flexible business model for the WiMAX technology, WiMAX forum identified and defined a series of business entities for the components of the WiMAX architecture that may, or may not, be implemented by the same real company. The defined business entities involved in the roaming process are [75]:

- Network Service Providers (NSPs) are business entities that provide IP connectivity and WiMAX services to WiMAX subscribers.
- Network Access Providers (NAPs) are business entities that provide WiMAX radio access infrastructure to one or more NSPs. NSPs may also have contractual agreements with other providers such as Internet Service Providers (ISPs).
- Home Network Service Provider (HNSP) is the service provider that has its users accessing the services of other operator's network through a roaming agreement.
- Visited Network Service Provider (VNSP) is the service provider that is hosting a node from another operator's network and with whom the VNSP has a roaming agreement.
- WiMAX Roaming Exchange (WRX) is an intermediary entity that can interconnect two or more NSPs to provide roaming service. NSPs may use the services of a WRX to handle specific functions while maintaining a bilateral roaming relationship with other NSPs, Hub Providers or Aggregators.

To enable a more broad and independent roaming process among operators the WiMAX forum defined WiMAX Roaming Interface (WRI). The definition of such interface does not prevent operators to exchange roaming information through proprietary interfaces, but it is a way to guarantee interconnection among different pairs that implements the interface. More details about the roaming process and the business and technical models defined by the WiMAX forum, to increase the coverage of NSPs, can be found in our work: [76].

---

## A.6 Conclusion

This appendix has addressed some of the most important aspects and challenges related to mobility management in WiMAX networks. The crucial concept for mobility management is the handoff, which is the process of transferring an ongoing session from one base station to another. The handoff has been studied in this appendix in all its forms: intra-WiMAX technology (horizontal handoff), inter-technologies (vertical handoff) and inter-providers (roaming). First we have described the different handoff mechanisms proposed by the IEEE 802.16e standard. Then, we have presented some of the works aiming at optimizing these procedures. We have classified the proposed works into two categories: those improving the handoff at layer 2 and those adopting an L2-L3 cross-layer approach in which the two layers collaborate to enhance the handoff performances. Among these cross-layer mechanisms, we have described more in details the fast MIPv6 handover mechanisms over 802.16e, proposed by IETF in [35].

The vertical handoff in heterogeneous networks—including WiMAX systems—has been considered first through some works proposed in the literature, and then through the MIH framework proposed by the IEEE 802.21 standard. Roaming, which is a key concept to increase the coverage of WiMAX network has also been briefly addressed in this appendix.

---



---

## Bibliography

- [1] IEEE Std 802.16-2004. IEEE Standard for Local and metropolitan area networks- Part 16: Air Interface for Fixed Broadband Wireless Access Systems. 2004.
  - [2] IEEE Std 802.16e 2005. IEEE Standard for Local and metropolitan area networks Part 16: Air Interface for Fixed and Mobile BWA Systems-Amendment 2: Physical and Medium Access Control Layers for Combined Fixed and Mobile Operation in Licensed Bands and Corrigendum 1. 2005.
  - [3] C. Cicconetti, A. Erta, L. Lenzini, and E. Mingozzi. Performance Evaluation of the IEEE 802.16 MAC for QoS Support. *IEEE Transactions on Mobile Computing*, 6(1):26–38, Jan. 2007.
  - [4] C. Cicconetti, L. Lenzini, E. Mingozzi, and C. Eklund. Quality of service support in IEEE 802.16 networks. *IEEE Network*, 20(2):50–55, Mar.-Apr. 2006.
  - [5] Alexander Sayenko, Olli Alanen, Juha Karhula, and Timo Hämäläinen. Ensuring the QoS requirements in 802.16 Scheduling. In *Proceedings of the 9th ACM international symposium on Modeling analysis and simulation of wireless and mobile systems*, pages 108–117, Terromolinos, Spain, 2006. ACM Press New York, NY, USA.
  - [6] Lin-Fong Chan, Hsi-Lu Chao, and Zi-Tsan Chou. Two-Tier Scheduling Algorithm for Uplink Transmissions in IEEE 802.16 Broadband Wireless Access Systems. In *International Conference on Wireless Communications, Networking and Mobile Computing, 2006. WiCOM 2006*, pages 1–4, Sept. 2006.
  - [7] Jianfeng Chen, Wenhua Jiao, and Hongxi Wang. A Service Flow management Strategy for IEEE 802.16 Broadband Wireless Access Systems in TDD Mode. In *IEEE International Conference on Communications (ICC2005)*, Seoul, Korea, 2005.
  - [8] Naian Liu, Xiaohui Li, Changxing Pei, and Bo Yang. Delay Character of a Novel Architecture for IEEE 802.16 Systems. In *Sixth International Conference on Parallel and Distributed Computing, Applications and Technologies, 2005. PDCAT 2005*, pages 293–296, Dec. 2005.
  - [9] Maode Ma, Jinchang Lu, S.K. Bose, and Boon Chong Ng. A three-tier framework and scheduling to support QoS service in WiMAX. In *6th International Conference on Information, Communications & Signal Processing*, pages 1–5, Dec. 2007.
  - [10] R. Perumalraja, J.J.J. Roy, and S. Radha. Multimedia Supported Uplink Scheduling for IEEE 802.16d OFDMA Network. In *Annual India Conference, 2006*, pages 1–5, Sept. 2006.
  - [11] M. Settembre, M. Puleri, S. Garritano, P. Testa, R. Albanese, M. Mancini, and V. Lo Curto. Performance analysis of an efficient packet-based IEEE 802.16 MAC supporting adaptive
-

- modulation and coding. In *International Symposium on Computer Networks, 2006*, pages 11–16, Jun. 2006.
- [12] J. Sun, Yanling Yao, and Hongfei Zhu. Quality of Service Scheduling for 802.16 Broadband Wireless Access Systems. In *IEEE 63rd Vehicular Technology Conference, 2006. VTC 2006-Spring*, volume 3, pages 1221–1225, 2006.
- [13] Kitti Wongthavarawat and Aura Ganz. Packet scheduling for QoS support in IEEE 802.16 broadband wireless access systems. *International Journal of Communication Systems*, 16(1):81–96, Feb. 2003.
- [14] Kitti Wongthavarawat and Aura Ganz. IEEE 802.16 based last mile broadband wireless military networks with quality of service support. *IEEE Military Communications Conference, 2003. MILCOM 2003*, 2(1):779–784, Oct. 2003.
- [15] N. Nasser and H. Hassanein. Prioritized multi-class adaptive framework for multimedia wireless networks. In *IEEE International Conference on Communications, 2004*, volume 7, pages 4295–4300, Jun. 2004.
- [16] D. Niyato and E. Hossain. Connection admission control algorithms for OFDM wireless networks. In *IEEE Global Telecommunications Conference, 2005. GLOBECOM '05*, volume 5, page 5 pp, Nov. 2005.
- [17] D. Niyato and E. Hossain. Delay-Based Admission Control Using Fuzzy Logic for OFDMA Broadband Wireless Networks. In *IEEE International Conference on Communications*, volume 12, pages 5511–5516, Jun. 2006.
- [18] Dusit Niyato and Ekram Hossain. Queue-Aware Uplink Bandwidth Allocation for Polling Services in 802.16 Broadband Wireless Networks. In *IEEE GLOBECOM 2005 proceedings*, 2005.
- [19] Dusit Niyato and Ekram Hossain. A Game-Theoretic Approach to Bandwidth Allocation and Admission Control for Polling Services in IEEE 802.16 Broadband Wireless Networks. In *3rd international conference on Quality of service in heterogeneous wired/wireless networks*, 2006.
- [20] Dusit Niyato and Ekram Hossain. Queue-Aware Uplink Bandwidth Allocation and Rate Control for Polling Service in IEEE 802.16 Broadband Wireless Networks. *IEEE TRANSACTIONS ON MOBILE COMPUTING*, 5(6), Jun. 2006.
- [21] Dusit Niyato and Ekram Hossain. A Queuing-Theoretic and Optimization-Based Model for Radio Resource Management in IEEE 802.16 Broadband Wireless Networks. *IEEE TRANSACTIONS ON COMPUTERS*, 55(11), Nov. 2006.
- [22] D. Niyato and E. Hossain. Joint Bandwidth Allocation and Connection Admission Control for Polling Services in IEEE 802.16 Broadband Wireless Networks. In *IEEE International Conference on Communications (ICC '06)*, volume 12, pages 5540–5545, Jun. 2006.
- [23] Vandana Singh and Vinod Sharma. Efficient and Fair Scheduling of Uplink and Downlink in IEEE 802.16 OFDMA Networks. In *IEEE Wireless Communications and Networking Conference, 2006. WCNC 2006*, volume 2, pages 984–990, Apr. 2006.
-

- 
- [24] Qingwen Liu, Xin Wang, and G.B. Giannakis. Cross-layer scheduler design with QoS support for wireless access networks. In *Second International Conference on Quality of Service in Heterogeneous Wired/Wireless Networks, 2005*, page 8 pp., Aug. 2005.
- [25] Qingwen Liu, Xin Wang, and G.B. Giannakis. A cross-layer scheduling algorithm with QoS support in wireless networks. *IEEE Transactions on Vehicular Technology*, 55(3):839–847, May 2006.
- [26] Yi-Ting Mai, Chun-Chuan Yang, and Yu-Hsuan Lin. Cross-Layer QoS Framework in the IEEE 802.16 Network. In *The 9th International Conference on Advanced Communication Technology*, volume 3, pages 2090–2095, Feb. 2007.
- [27] Yi-Ting Mai, Chun-Chuan Yang, and Yu-Hsuan Lin. Design of the Cross-Layer QoS Framework for the IEEE 802.16 PMP Networks. *IEICE Transactions on Communications*, E91-B(5):1360–1369, May 2008.
- [28] K. A. Noordin and G. Markarian. Cross-Layer Optimization Architecture for WiMAX Systems. In *The 18th Annual IEEE International Symposium on Personal, Indoor and Mobile Radio Communications (PIMRC'07)*, pages 1–4, Sept. 2007.
- [29] Dionysia-Katerina Triantafyllopoulou, Nikos Passas, and Alexandros Kaloxylos. A Cross-Layer Optimization Mechanism for Multimedia Traffic over IEEE 802.16 Networks. In *European Wireless 2007 (EW 2007)*, Apr. 2007.
- [30] Dionysia Triantafyllopoulou, Nikos Passas, Apostolis K. Salkintzis, and Alexandros Kaloxylos. A heuristic cross-layer mechanism for real-time traffic over IEEE 802.16 networks. *INTERNATIONAL JOURNAL OF NETWORK MANAGEMENT*, 17(5):347–361, Sept. 2007.
- [31] Scalable Network Technologies. Qualnet 4.5, March 2008. <http://www.scalable-networks.com/products/qualnet/>.
- [32] European Telecommunications Standards Institute (ETSI). European profile standard for the physical and medium access control layer of Intelligent Transport Systems operating in the 5 GHz frequency band, June 2010. ETSI ES 202 663 V1.1.0. Intelligent Transport Systems (ITS).
- [33] R.Q. Hu, D. Paranchych, Mo-Han Fong, and Geng Wu. On the evolution of handoff management and network architecture in WiMAX. In *Proceedings of the Mobile WiMAX Symposium, 2007*, pages 144–149, Mar. 2007.
- [34] WiMAX Forum. WiMAX Forum Network Architecture (Stage 2: Architecture Tenets, Reference Model and Reference Points) [Part 1]. volume Release 1, Version 1.2, Jan. 2007.
- [35] H. Jang, J. Jee, Y. Han, S. Park, and J. Cha. Mobile IPv6 Fast Handovers over IEEE 802.16 Networks. RFC 5270 (Informational), June 2008.
- [36] IEEE 802.21-2008. IEEE Standard for Local and metropolitan area networks Part 21: Media Independent Handover. Jan. 2008.
- [37] IEEE 802.16-2009. IEEE Standard for Local and metropolitan area networks Part 16: Air Interface for Broadband Wireless Access Systems. May 2009.
-

- 
- [38] B. Gi Lee and S. Choi. *Broadband Wireless Access And Local Networks: Mobile WiMAX and WiFi*. Artech House Publishers, 2008.
- [39] WiMAX Forum. WiMAX Forum Mobile System Profile Specification - FDD Specific Part . In *WMF-T23-003-R015v01*, volume Release 1.5, Version 1, Jan. 2009.
- [40] IEEE Std 802.16a 2003. IEEE Standard for Local and metropolitan area networks Part 16: Air Interface for Fixed Broadband Wireless Access Systems. 2003.
- [41] D. Stiliadis and A. Varma. Latency-rate servers: a general model for analysis of traffic scheduling algorithms. *IEEE/ACM Transactions on Networking*, 6(5):611–624, Oct. 1998.
- [42] H. Fattah and C. Leung. An overview of scheduling algorithms in wireless multimedia networks. *IEEE Wireless Communications*, 9(5):76–83, Oct. 2002.
- [43] Howon Lee, Taesoo Kwon, and Dong-Ho Cho. An Enhanced Uplink Scheduling Algorithm Based on Voice Activity for VoIP Services in IEEE 802.16d/e System. *IEEE COMMUNICATIONS LETTERS*, 9(8), Aug. 2005.
- [44] Howon Lee, Taesoo Kwon, Dong-Ho Cho, Geunhwi Limt, and Yong Changt. Performance Analysis of Scheduling Algorithms for VoIP Services in IEEE 802.16e Systems. In *IEEE 63rd Vehicular Technology Conference, 2006. VTC 2006-Spring*, volume 3, pages 1231–1235, 2006.
- [45] Yin Ge and Geng-Sheng Kuo. An Efficient Admission Control Scheme for Adaptive Multimedia Services in IEEE 802.16e Networks. In *IEEE 64th Vehicular Technology Conference, 2006. VTC-2006 Fall. 2006*, pages 1–5, Sept. 2006.
- [46] D. Niyato and E. Hossain. Radio Resource Management Games in Wireless Networks: An Approach to Bandwidth Allocation and Admission Control for Polling Service in IEEE 802.16. *IEEE Wireless Communications*, 14(1), Feb 2007.
- [47] H. Wang, W. Li, and D. P. Agrawal. Dynamic admission control and QoS for 802.16 wireless MAN. In *Wireless Telecommunications Symposium, 2005*, pages 60–66, Apr 2005.
- [48] Liping Wang, Fuqiang Liu, Yusheng Ji, and Nararat Ruangchaijatupon. Admission Control for Non-preprovisioned Service Flow in Wireless Metropolitan Area Networks. In *Fourth European Conference on Universal Multiservice Networks, 2007. ECUMN '07*, pages 243–249, Feb. 2007.
- [49] Chi-Hong Jiang and Tzu-Chieh Tsai. Token bucket based CAC and packet scheduling for IEEE 802.16 broadband wireless access networks. In *3rd IEEE Consumer Communications and Networking Conference, 2006. CCNC 2006*, volume 1, pages 183–87, Jan. 2006.
- [50] S. Chandra and A. Sahoo. An Efficient Call Admission Control for IEEE 802.16 Networks. In *Proceedings of the 15th IEEE LAN/MAN Workshop, LANMAN 2007*, pages 188–193, Jun. 2007.
- [51] E. Kwon, J. Lee, K. Jung, and S. Ryu. A Performance Model for Admission Control in IEEE 802.16. In *3rd International Conferences on Wireless/Wired Internet Communications (WWIC 2005)*, pages 159–168, Xanthi , Greece, May 2005.
-

- 
- [52] J.Y. Lee and K.B. Kim. Statistical Admission Control for Mobile WiMAX Systems. In *IEEE Wireless Communications and Networking Conference. WCNC 2008*, Apr. 2008.
- [53] O. Yang and J. Lu. A New Scheduling and CAC Scheme for Real-Time Video Application in Fixed Wireless Networks. In *3rd IEEE Consumer Communications and Networking Conference, 2006. CCNC 2006*, volume 1, pages 303–307, Jan 2006.
- [54] O. Yang and J. Lu. Call Admission Control and Scheduling Schemes with QoS Support for Real-time Video Applications in IEEE 802.16 Networks. *JOURNAL OF MULTIMEDIA (JMM)*, 1(2):21–29, May 2006.
- [55] Jaiyong Lee Eunhyun Kwon, Hun-je Yeon and Kyunghun Jung. Markov Model for Admission Control in the Wireless AMC Networks. *IEICE Transactions on Communications*, E89-B(8):2230–2233, 2006.
- [56] W. Fang, N. Seddigh, and B. Nandy. A Time Sliding Window Three Colour Marker (TSWTCM). RFC 2859 (Experimental), June 2000.
- [57] Richard Bossom et al. D31 European ITS Communication Architecture - Overall Framework - Proof of Concept Implementation, March 2009. COMeSafety deliverable.
- [58] IEEE P802.11p/D7.0. IEEE Draft Standard for Information Technology – Telecommunications and information exchange between systems – Local and metropolitan area networks – Specific requirements, Part 11: Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) specifications, Amendment 7: Wireless Access in Vehicular Environments. May 2009.
- [59] ISO TC204 Working Group 16. Continuous Air interface for Long and Medium range (CALM), Sept. 2009. <http://www.isotc204wg16.org/wg16>.
- [60] European Telecommunications Standards Institute (ETSI). <http://www.etsi.org/WebSite/technologies/IntelligentTransportSystems.aspx>, 2010.
- [61] Centre for Applied Informatics (ZAIK) and the Institute of Transport Research at the German Aerospace Centre. Simulation of Urban MObility (SUMO) 0.9.8, Feb. 2008. <http://sumo.sourceforge.net/>.
- [62] Feliz Kristianto Karnadi, Zhi Hai Mo, and Kun chan Lan. Rapid Generation of Realistic Mobility Models for VANET. In *In Proc. of the IEEE Wireless Communications and Networking Conference, 2007. WCNC 2007*, March 2007. <http://sumo.sourceforge.net/>.
- [63] Doo Hwan Lee, K. Kyamakya, and J.P. Umondi. Fast handover algorithm for IEEE 802.16e broadband wireless access system. In *Proceedings of the 1st International Symposium on Wireless Pervasive Computing, 2006*, page 6 pp., Jan. 2006.
- [64] Jenhui Chen, Chih-Chieh Wang, and Jiann-Der Lee. Pre-Coordination Mechanism for Fast Handover in WiMAX Networks. In *Proceedings of the The 2nd International Conference on Wireless Broadband and Ultra Wideband Communications, 2007. AusWireless 2007*, Aug. 2007.
- [65] Ling Chen, Xuejun Cai, Rute Sofia, and Zhen Huang. A Cross-Layer Fast Handover Scheme For Mobile WiMAX. In *Proceedings of the IEEE 66th Vehicular Technology Conference, 2007. VTC-2007 Fall*, pages 1578–1582, Oct. 2007.
-

- 
- [66] Chung-Kuo Chang and Chin-Tser Huang. Fast and Secure Mobility for IEEE 802.16e Broadband Wireless Networks. In *Proceedings of the International Conference on Parallel Processing Workshops, 2007. ICPPW 2007*, pages 46–52, Sept. 2007.
- [67] Ed. R. Koodli. Mobile IPv6 Fast Handovers. RFC 5268 (Standards Track), 2008.
- [68] Mohamed KASSAB. *Layer-2 Handover Optimization for Intra-technologies and Inter-technologies Mobility*. PhD thesis, Institut Télécom-Télécom Bretagne, UR1 - Université de Rennes 1, 2008.
- [69] Youngkyu Choi and Sunghyun Choi. Service Charge and Energy-Aware Vertical Handoff in Integrated IEEE 802.16e/802.11 Networks. In *Proceedings of the 26th IEEE International Conference on Computer Communications. IEEE INFOCOM 2007*, pages 589–597, May 2007.
- [70] Janise McNair and Fang Zhu. Vertical handoffs in fourth-generation multinet network environments. *IEEE Wireless Communications*, 11(3):8–15, June 2004.
- [71] Z. Daia, R. Fracchiaa, J. Gosteaub, P. Pellatia, and G. Vivier. Vertical handover criteria and algorithm in IEEE 802.11 and 802.16 hybrid networks. In *Proceedings of the IEEE International Conference on Communications, 2008. ICC'08*, pages 2480–2484, May 2008.
- [72] Yu Liu and Chi Zhou. A Vertical Handoff Decision Algorithm (VHDA) and a Call Admission Control (CAC) policy in integrated network between WiMax and UMTS. In *Proceedings of the Second International Conference on Communications and Networking in China, 2007. CHINACOM'07*, pages 1063–1068, Aug. 2007.
- [73] Lambros Sarakis, George Kormentzas, and Francisco Moya Guirao. Seamless Service Provision For Multi Heterogeneous Access. *IEEE Wireless Communications*, Oct. 2009.
- [74] E. Piri and K. Pentikousis. IEEE 802.21: Media Independent Handover Services. *The Internet Protocol Journal*, 12:7–27, June 2009.
- [75] WiMAX Forum. WiMAX Forum Roaming Models White Paper. In *WMF-T48-001-v01 Approved Version 1*, volume Version 1, Apr. 2009.
- [76] Ikbal Chammakhi Msadaa, Daniel Câmara, and Fethi Filali. *WiMAX Security and Quality of Service : An End-to-End Perspective*, chapter Mobility Management in WiMAX Networks. WILEY, 2010.
-