

# On performance of TCP and VoIP traffic in mobile WiMAX networks

DRAGORAD MILOVANOVIC, ZORAN BOJKOVIC

University of Belgrade,  
Studentski trg 1, 11000 Belgrade, Serbia  
e-mail: dragoam@gmail.com, <http://dragorad.milovanovic.googlepages.com>

*Abstract:* - The major performance objectives of the next-generation wireless communication systems are high link capacity and increase application/service performance from the user perspective. An attractive technology for providing broadband access is WiMAX network. For the network and application-level capacity and performance analysis, we first provide an overview of WiMAX network architecture and system performance. Afterward we discuss the benefits and challenges of flat architecture for mobile networks. In this kind of architecture the access service network gateway and base stations are consolidated into a single channel. All-IP flat architecture is a promising option. However, performance goals of the WiMAX system are conflicting: maintenance VoIP quality vs. completely utilization of the remaining capacity for TCP. We investigate the approaches in performance optimization based on TCP modification schemes. The solution is to control TCP traffic before entering WiMAX network in conjunction with methods to dynamically estimate available bandwidth in real time.

Key-Words: - Wireless networking, TCP performance, Traffic control.

## 1 Introduction

The success of WiMAX high data rate communications in metropolitan area networks (MAN) depends on its capability of providing cost-effective solution for a variety of services [1]. In 2001, the first IEEE802.16 standard was published, while in 2005, the standard IEEE802.16e was approved as the official standard for mobile applications [2]. The mobile WiMAX systems have a higher system capacity and a more sophisticated mechanism to provide a better quality of service (QoS) [3]. To evaluate the mobile WiMAX system capacity and performance, all the aspects of the performance evaluation – from air link to application – are required.

A mobile WiMAX network support both voice and TCP services. However, performance goals of the system are conflicting: maintenance VoIP quality vs. completely utilization of the remaining capacity for TCP. Services VoIP and TCP cannot simply share the WiMAX medium without severe voice quality degradation and/or reduction in TCP capacity [4]. In order to investigate the interaction between these two categories of traffic, we present TCP modification and control schemes. The solution is to control TCP traffic before entering the multihop network.

The article is structured as follows. The next section discusses the state of the art in WiMAX networks and system performance evaluation, together with benefits of flat architecture of mobile network. We then present solution approaches for TCP and VoIP performance optimization based on traffic modification and control.

## 2 WiMAX networks architecture

WiMAX (*Worldwide interoperability for microwave access*) has been proposed as an attractive wireless communication technology due to the fact that it can provide high data rate communications for metropolitan areas. Until now, a number of specifications for WiMAX were standardized by the IEEE802.16 Working group [5]. In addition companies in the industry also have formed the *WiMAX Forum* to promote the development and deployment of WiMAX systems. According to the standards, WiMAX can support up to a 75 Mbps data rate (single channel) and can cover on how it can provide cost-effective solutions for a variety of existing and potential services [6].

The specifications in the current WiMAX standards can be partitioned into two important parts, the physical (PHY) layer and the medium access control (MAC) layer:

- IEEE 802.16 supports four PHY specifications for the licensed bands: wireless MAN-SC (Single Carrier, 10-66 GHz), Wireless MAN-SCa (Single Carrier below 11 GHz), Wireless MAN-OFDM (Orthogonal Frequency Division Multiplexing, below 11 GHz), Wireless MAN-OFDMA (Orthogonal Frequency Division Multiple Access, below 11 GHz).
- In the MAC layer, IEEE 802.16 supports two modes: the point-to-multipoint (PMP) mode and the mesh mode. PMP network is designed primarily for providing the last-mile access to the Internet service providers (ISPs). It can be viewed as a tree topology in which the root is a base station (BS), while the subscriber stations (SSs) are the leaves. As an example of this network topology, in Figure 1, we have WiMAX network architecture with PMP mode and mesh mode. The PMP mode includes one base station and a number of subscriber stations.

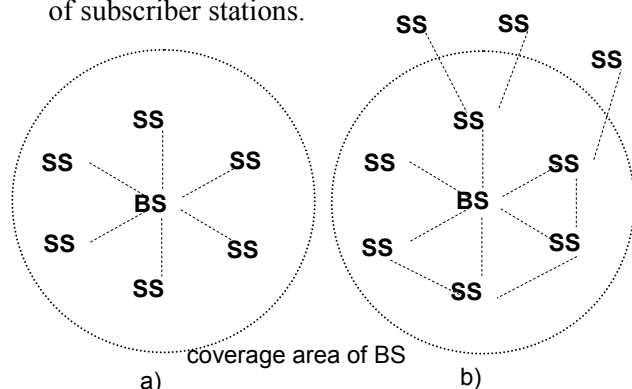


Figure 1. WiMAX network architecture: a) PMP mode, b) mesh mode.

In the MAC layer, IEEE802.16 supports two modes, the point-to-multipoint (PMP) mode and the mesh mode:

- PMP network is designed primarily for providing the last-mile access to the Internet service providers (ISPs). It can be viewed as a tree topology in which the root is a base station (BS), while the subscriber stations (SSs) are the leaves. As an example of this network topology, in Figure 1, we have WiMAX network architecture with PMP mode and mesh mode. The PMP mode includes one base station and a number of subscriber stations.
- The WiMAX mesh network is a multihop *ad hoc* network in which subscriber stations can connect with one another directly. In the case of mesh mode, BS can provide access to the service provider. A relay station (RS) is a special type of SS in that it can forward traffic flows to BSs or other SSs and a mobile station (MS) is an SS that

can move within network. The concept of MS is specified in IEEE802.16e, which extends the PMP mode and defines a concept of mobile multihop relay (MMR) networking. One reason for proposing the MMR scheme is to overcome the problem that the mesh mode is not compatible with the PMP mode. Comparing to the PMP mode, the mesh mode is more flexible and can be used to deploy the infrastructure.

The characteristics of major services that are crucial for WiMAX networks are compared in Table 1. The services do not comprise a complete list. Other services, such as online gaming, telemedicine/e-health services are also important. Also, security and reliability generally should be required by all services.

Services	data rate	delay	content access	commun. pattern
Internet access	fixed (for maximum)	best effort	no	one-to-one
VoIP	fixed or variable	real-time	no	one-to-one and many-to-many
Videoconferencing	variable	real-time	no	many-to-many
Media downloading	variable	nonreal-time	yes	one-to-one and many-to-many
IPTV	fixed or variable	nonreal-time	yes	one-to-one and many-to-many

Table 1. Characteristics of major services for WiMAX networks.

To meet the requirements of existing and potential services in WiMAX networks the design guidelines like connection-orientation, group communications, security, reliability and storage have to be proposed.

WiMAX networks must provide group-based communications efficiently. Many important services will require group-based communications, including one-to-many and many-to-many modes. Because WiMAX networks are likely to be multi-hop wireless networks, it is important to address the security issues in the network layer, in addition to the MAC layer. Nevertheless, the solution can be based on just the network layer, and it can be performed in a cross-layer manner. Also, due to the nature of wireless communications, it is important to provide reliability in the design of WiMAX networks. Finally, nodes in WiMAX networks can provide storage capacity. In this way, delay-insensitive content distribution may be supported effectively in WiMAX networks.

Currently deployed mobile WiMAX networks use hierarchical architecture. Flat architecture is specified as a design alternative in the mobile WiMAX standard. In flat architecture, the functionalities of

the controller and BS are consolidated into a single element. The integrated element is directly connected to the IP (Internet Protocol) core networks. This is better suited to deal with bursty Internet traffic than the traditional hierarchical architecture. The main benefits of flat architecture for providing wireless data services can be summarized as follows:

- The controller – integrated BSs are connected directly to the core network (Internet), eliminating the requirements of the wireless technology – dependent radio-access network (RAN). This is the access service network (ASN) in mobile WiMAX networks between the BS and the IP core network. The interoperability with heterogeneous wireless technologies is easy to achieve because the BSs are directly connected to the wireless technology-neutral IP core network.
- Flat architecture provides highly scalability because there is new centralized performance bottleneck, while traffic is processed in a fully distributed fashion. It also provides adding or removing cells without concern about the capacity of the centralized controller to which the new BSs are attached.
- Integrated design reduces the cost of intermodulate communication (sophisticated cross-layer optimization is possible for performance gain).
- Flat architecture achieves resource efficiency gains by preventing suboptimal routing. Under the hierarchical architecture, all traffic must pass through the centralized controllers, which may extend the routing path, resulting in suboptimal traffic routing.
- Under flat architecture, single failure points (the centralized controller) do not exist, and the impact of a BS failure can be limited locally.
- Flat architecture has economic advantages. The general purpose IP equipment is much cheaper than radio access network components because of economics of scale. The options for network management tools for IP networks are available at low cost as well.

The challenges that the flat architecture networks face are: handover performance, quality of service (QoS), self-configuration and self-optimization, self-diagnosis. A high performance IP-mobility solution for flat architecture is needed. In flat architecture, network configuration and resource management should be done in a distributed way, namely, self-configuration and self-optimization. This challenge comes with the benefits of scalability and flexibility.

## 3 WiMAX system performance

### 3.1 Performance requirements

Mobile WiMAX introduces OFDMA and supports several key features necessary for delivering mobile broadband services at vehicular speeds greater than 120 kmh with QoS comparable to broadband wireline access alternatives. These features and attributes include:

- tolerance to Multipath and Self-Interference with subchannel orthogonality in both the DL and the UL
- scalable Channel Bandwidths from 1.25 to 20 MHz
- Time Division Duplex (TDD) is defined for the initial mobile WiMAX profiles for its added efficiency in support of asymmetric traffic and channel reciprocity for easy support of advanced antenna systems.
- Hybrid-Automatic Repeat Request (H-ARQ) provides added robustness with rapidly changing path conditions in high mobility situations.
- Frequency Selective Scheduling and subchannelization with multiple permutation options, gives mobile WiMAX the ability to optimize connection quality based on relative signal strengths to specific users.
- Power Conservation Management ensures power-efficient operation of battery operated mobile handheld and portable devices in Sleep and Idle modes.
- Network-Optimized Hard Handoff (HHO) is supported to minimize overhead and achieve a handoff delay of less than 50 milliseconds.
- Smart Antenna support aided by subchannelization and channel reciprocity enables a wide range of advanced antenna systems including beamforming, space-time coding and spatial multiplexing
- Fractional Frequency Reuse controls co-channel interference to support universal frequency reuse with minimal degradation in spectral efficiency.
- 5 millisecond Frame Size provides optimal tradeoff between overhead and latency.

The following types of mobile environments are considered:

- Nomadic. Where the subscriber station may periodically connect from different points in the network and may involve different base stations. Handover support for a live connection is not required in this case.
- Portable. The subscriber station is a portable device, and handover support is required as it moves across base stations.
- Simple mobility. This is a mode supports mobile stations with speeds up to 60kmh and a hard handover with possible brief interruption.
- Full mobility. Support of mobility at vehicular speeds (up to 120kmh) and a seamless handover.

The mobile environment presents many challenges:

- Signals received at mobile devices can vary greatly over short distances and can contain severe errors. Mobile WiMAX uses OFDM/OFDMA modes of transmission, beamforming and MIMO to have spatial diversity, better link margins, and reduction of interference.
- Mobile devices are constrained by battery power and processing capabilities. Fixed WiMAX common TDM downlink frame for all active devices is changed to a TDMA mode with downlink bursts meant for individual mobile devices or device groups.
- In Fixed WiMAX, the subscriber devices are assigned a fixed time slot for uplink transmission, at which time they must transmit on all subcarriers simultaneously. In Mobile WiMAX, the mobile device needs to transmit only on the assigned subcarriers, reducing the peak power requirements.
- The mobile devices by their very nature can move in and out of various areas which may be operating fixed services in unlicensed bands. The mobile WiMAX has been specified only in specific licensed bands.

Mobile WiMAX adds significant enhancements:

- improved NLOS coverage by utilizing advanced antenna diversity schemes and hybrid automatic repeat request (HARQ),
- adopted dense subchannelization, thus increasing system gain and improving indoor penetration,
- used adaptive antenna system (AAS) and multiple input multiple output (MIMO) technologies to improve coverage,
- introduced a downlink subchannelization scheme, enabling better coverage and capacity trade-off.

### 3.2 System performance

System **capacity** is an important factor affecting the performance of a WiMAX networks. Capacity determines the amount of data that can be delivered to and from the users. There is a limit on how much data can be reliably sent through any given channel. There are several ways of quantifying the capacity of a wireless system:

- the traditional way of quantifying capacity is by calculating the data rate per unit bandwidth that can be delivered in a system
- an alternate way of looking at capacity is by quantifying the number of users that can be supported by a system, by a sector, or per channel to name a few.

A channel can also be defined in several ways: it can be defined as a frequency range, a time slot, or as a frequency-time combo slot.

System capacity is dependent on the distribution of the users in the service area as well as the type of

service that is requested. In systems that employ adaptive modulation such as WiMAX, capacity is also a function of the C/I ratio since higher modulation orders can achieve higher spectral efficiencies. Channel and, if necessary, sub-channel planning also has an effect on the data rate that can be supported per unit bandwidth. The flexibility in the WiMAX standard allows a system designer to perform tradeoffs between overall system capacity and interference levels to best allocate the resources to supply the customers with their data needs.

Capacity prediction of WiMAX system requires classification of users based on their demands and how much load they place on the system. Browsing the web, emailing, sending/receiving video, downloading files, or using VoIP are all applications that might be performed simultaneously within the population of users. Each of these applications places varying demands on the system. Some of them might require a higher data rate on download than on upload, while others are about evenly distributed. At any given time there will be a mixture of services that are being requested and the system is able to support several voice users with the same resources it takes to serve one video user.

An alternate way of looking at capacity is to establish the load that a typical number of users place on a system and then determine at what point the load surpasses each sector's ability to deliver. The number of system resources that are consumed in a given area depends on its demographics and the type of users that are present in that location (factors such as terrain and time of day can also affect demand). Additionally, in the case of WiMAX, the signal levels received at the mobile and base station terminals are important since users that achieve better C/I ratios can be reached using higher order modulation schemes therefore consuming less of available slots in a given sector.

It can be useful to examine capacity in the extreme cases where a sector can achieve the maximum or minimum of both the number of supported users and spectral efficiency (bits/second/Hz) in a sector. The worst case occurs when all the users are at the edge of the coverage area of the cell and are only reachable by the lowest order modulation and at the same time have the greatest demand on the system. On the other hand, the most number of supported users and the best spectral efficiencies occur when all users are close to the base station and demand low data rate services.

An analytical framework that allows to evaluate the amount of supported VoIP users as well as throughput achievable with TCP/IP connections as a function of the most significant parameters characterizing WiMAX technology is given in [7].

In order to evaluate the maximum **number of users** performing a VoIP call that can be served by WiMAX in a function of the adopted codec, the scheduling service and the transmission mode, all users are supposed to be served adopting the same mode. Also in this case we need to consider the whole packet processing from the application layer down to the transmission over the medium. Before being transmitted, each packet generated by the codec must pass through the whole protocol stack; the RTP, UDP, IP and MAC layers overheads are added. Next, each packet to be transmitted is mapped onto OFDMA-slots, thus, in order to assess the maximum number of users that can be served, the number of slots that are needed for each packet must be firstly calculated. Depending on the considered scheduling service and the specific VoIP codec, the average number of slots required in a given (DL/UL) subframe by a single user adopting mode, is given by analytical expression for the considered scheduling services [7].

In order to evaluate the maximum WiMAX end-to-end system **throughput** at the application layer, a single TCP connection is supposed to be active either in the downlink or in the uplink. Furthermore, the return link (the uplink when considering a downlink TCP flow and vice versa) is considered to be always sufficient for the transmission of TCP acknowledgments; this assumption implies that the analysis is focused on the direction of the data flow. The amount of application layer data conveyed by a single MSDU (MAC service data unit) is given by  $S_{AS}$  bytes. It follows that the average amount of application layer data accommodated in a given subframe is:

$$S = N_{MS} \cdot S_{AS} \quad (1)$$

where  $N_{MS}$  is number of complete MSDUs that can be allocated in the considered subframe. The system application layer throughput per one TCP connection can be derived as:

$$TCP_{throughput} [bps] = \frac{8 \cdot S [bytes]}{T_F [s]} \quad (2)$$

where  $T_F$  represents the frame duration.

The main outcome of the analysis in [7] is given by a set of criteria to be followed in order to maximize the WiMAX throughput provided to the VoIP users.

## 4 VoIP and TCP traffic

### 4.1 Voice over wireless IP

VoIP integration of conventional telephone services with the growing number of other IP-based applications is seen as one of the important technology for telecommunications providers. In addition to the cost reduction achievable by the sharing of network resources, VoIP is accepted to accelerate the development of rich multimedia services. Since quality is not generally guaranteed in an IP network, it is important for the networks and terminals to be properly designed before providing services. The quality of service had to be constantly monitored, while action is necessary to maintain the level of service. Since VoIP systems are based on new coding technologies and a new transmission technology, the primary determinants of the perceptual QoS of VoIP service are distortions caused by speech coding and packet losses, loudness, delay and echo [8].

Wireless networks are providing high bandwidths to users, enabling real-time multimedia applications. Expected improvements with respect to current third-generation (3G) networks include data rates up to 100 Mbps in wide area networks (WANs) and 1 Gbps in wireless local area networks (WLANs). Multimedia applications may generate real-time high bit rate flows. Applications generating small payloads, such as VoIP, are responsible for the highest relative overheads. A typical IPv6 VoIP packets includes a 40-byte IP-header, an 8-byte UDP header, a 12-byte RTP header, and a 20-byte payload. In this case, the total IP overhead is about three times the size of the payload, or 75% of the total packet size.

Assuming that the voice traffic currently transported by systems such as the GSM (*Global System for Mobile Communication*) migrates to 4G packet switched networks, the amount of VoIP traffic is expected to become significant. Combining significant amounts of real-time traffic and large relative overheads, we can justify the study of header compression (HC) techniques and their value in future wireless networks.

Header compression techniques, such as robust header compression can be used to reduce the overhead of IP-based traffic. The robust HC (RoHC) standard describes a set of HC mechanisms for compressing the headers of IP-based protocols. RoHC mechanisms were initially defined for wireless links with high bit error rates (BERs). In opposite, new wireless data networks, such as IEEE802.11 offer a confirmed frame delivery service that reduce

the transmission failure to insignificant values, although at the expense it higher delays. In IEEE802.11 links, the access to the medium is regulated by a distributed algorithm named DCF (Distributed Coordination Function), which is based on contention. The use of short packets, such as those produced by VoIP, decreases the overall transmission efficiency of the 802.11 link.

#### 4.2 TCP performance in wireless networks

The transport protocol family for media streaming includes UDP (User Datagram Protocol), TCP (Transmission Control Protocol), RTP (Real-Time Protocol) and RTCP (Real-Time Control Protocol). UDP and TCP provide basic transport functions, while RTP and RTCP run on top of UDP/TCP. UDP and TCP protocols support such functions as multiplexing, error control, congestion control or flow control:

- UDP and TCP can multiplex data streams for different applications running on the same machine with the same IP address.
- For the purpose of error control, TCP and most UDP implementations employ the checksum to detect bit error. If single or multiple bit errors are detected in the incoming packet, the TCP/UDP layer discards the packet so that the upper layer (RTP) will not receive the corrupted packet. On the other hand, different from UDP, TCP uses retransmission to recover lost packets. Therefore, TCP provides reliable transmission, while UDP does not.
- TCP employs congestion control to avoid sending too much traffic, which may cause network congestion.
- TCP employs flow control to prevent the receiver buffer from overflowing, while UDP does not have any flow control mechanisms.

A large amount of research has focused on the optimization of TCP performance in wireless networks [4]. There are three main categories: connection splitting, link layer solutions, and gateway solutions.

- **Connection splitting** can hide the wireless link by terminating the TCP connection prior to the wireless at the base station or access point. In that way, the communication in a wireless network can be optimized independent of the TCP applications. However, it requires an extra overhead to maintain two connections for each flow. It also violates end-to-end TCP semantics and requires a complicated handover process.

- **Link layer solutions.** The idea is to make the wireless link layer look similar to the wired case from the perspective of TCP. The relevant and interesting proposal is the so called snoop protocol. Here, the snoop agent is introduced at the base station (BS) to perform local retransmission using information sniffed from the TCP traffic passing through the BS. Another link layer solution proposed QoS scheduling with priority queues within the intermediate nodes of a multihop network to improve VoIP quality by placing TCP data in a lower QoS level.
- **Gateway solutions.** One way to address TCP performance problems within wireless networks is to evenly space or pace data sent into the multihop over an entire roundtrip time so that data are not sent in a burst. Pacing is implemented using a data and/or ACK pacing mechanism.

#### 4.3 VoIP and TCP traffic interaction

The interaction between TCP and VoIP over wireless network is complex [4, 9]. First of all, TCP is an end-to-end protocol. There are no explicit signaling mechanisms in the network to tell TCP peers how fast to send, how much to send, or when to slow down a transmission. TCP needs to be aggressive in discovering link bandwidth because that is how it can achieve high utilization.

TCP produces burst traffic, while VoIP produces uniform (stream) traffic. When the network is congested by interference or too much TCP data, VoIP traffic suffers from increased network losses and delays. However, TCP goes into the receiving stage, reducing its sending rate until the network recovers from congestion, and then sends all postponed packets. This cycle of burstiness leads to both low utilization for TCP and unacceptable quality for voice.

TCP assumes that losses come from congestion. This observation has been the basis for many studies focusing on preventing the TCP congestion control mechanism reacting to link layer errors. Also, TCP behavior may lead to poor performance because of packet drops due to hidden terminal induced problems such as channel interference and TCP data/acknowledgement (ACK) convention.

VoIP packets are small, while, TCP packets are large. For a given bit error rate, TCP packets have less success, so many of them would be retransmitted across multihop links, thus generating even more load that in turn generates more interference. A wireless node cannot determine by itself what the interference conditions are in its neighborhood.

Factors affecting regional interference are actual paths, load, physical distance between nodes, collision domains, and hidden terminal.

VoIP is mostly constant bit rate, has tight delay and loss requirements, and should always be served prior to TCP traffic. Classical solutions such as priority queues, bandwidth limitation, and traffic shopping do not provide satisfactory solutions for the coexistence problem. If voice traffic has priority locally within a node, bursty TCP traffic affects voice packets on other nodes within the interference range.

Bursty traffic produces higher queuing delays, more packet loss and lower throughput. TCP congestion control mechanisms and self-clocking create extremely bursty traffic in networks with large bandwidth-delay products, cause long queues and increase the likelihood of massive losses. Mobile WiMAX network traffic tends to have self-similar behavior, which is harmful to traffic requiring a stable bit rate, such as VoIP or streaming.

As an example, the hybrid wire/wireless network is shown in Fig.2. The multihop extension between forwards TCP traffic from wired Internet and VoIP calls to/from an IP-private branch exchange (PBX) through the gateway. TCP data are flowing from the gateway G to the client. The multihop leg is where VoIP needs to be protected from TCP. As an access network for VoIP and TCP, WiMAX 802.11 based multihop is used. Downstream TCP traffic in which data flows from servers across the Internet, through the gateway to wireless clients is considered. On the other hand, TCP acknowledgments (ACK) travel in the opposite direction. VoIP traffic is symmetric and bidirectional. It is clear that VoIP and TCP cannot simply share a multihop network without experiencing severe voice quality degradation and/or reduction in TCP capacity.

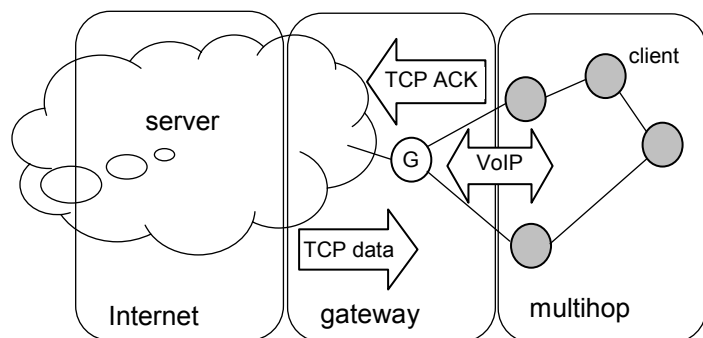


Figure 2. Access network for VoIP and TCP.

## 5 TCP modification and control schemes

The efficient support of IP communications in wireless environments is considered a key issue of WiMAX multimedia networks. The IP protocol and its main transport layer companions (TCP and UDP) were also designed for fixed networks, with the assumption that the network consists of point-to-point physical links with stable available capacity. However, when a wireless access technology is used in the link layer, it could introduce severe variations on available capacity, and could thus result in low TCP protocol performance [1]. There are two main weaknesses of the IP over wireless links:

- The assumption of reliable communication links.** Assuming highly reliable links, the only cause of undelivered IP packets is congestion at some intermediate nodes, which should be treated in higher layers with an appropriate end-to-end congestion control mechanism. UDP, targeted mainly for real-time traffic, does not include any congestion control, as this would introduce unacceptable delays. Instead, it simply provides direct access to IP, leaving applications to deal with the limitations of IP's best effort delivery service. TCP, on the other hand, dynamically tracks the round-trip delay on the end-to-end path and times out when acknowledgments are not received in time, retransmitting unacknowledged data. Additionally, it reduces the sending rate to a minimum and then gradually increases it in order to probe the network's capacity. In WLANs, where errors can occur due to temporary channel quality degradation, both these actions (TCP retransmissions and rate reduction) can lead to increased delays and low utilization of the source available bandwidth.
- The lack of traffic prioritization.** Designed as a 'best effort' protocol, IP does not differentiate treatment according to the kind of traffic. For example, delay sensitive real-time traffic, such as VoIP, will be treated in the same way as ftp or e-mail traffic, leading to unreliable service. In fixed networks, this problem can be relaxed with over-provisioning of bandwidth, wherever possible (e.g., by introducing high capacity fiber optics). In WLANs this is not possible because the available bandwidth can be as high as a few tens of Mbps. But even if bandwidth was sufficient, multiple access could still cause unpredictable delays for real-time traffic. For these reasons, the introduction of scheduling mechanisms is required for IP over WLANs, in order to ensure reliable service under all kinds of conditions.



### 5.1 TCP enhanced variants

We first consider the selected TCP variants in typical downstream traffic scenario in which TCP sources across the Internet send bulk traffic to a user connected to network. Basically, TCP is a reliable, connection-oriented data transfer protocol, with well-established mechanisms for flow and congestion control based on the sliding window [10].

**TCP Reno** is classic TCP scheme. It has four transmission phases: slow start, congestion avoidance, fast recovery, and fast retransmit. Reno employs a sliding-window-based flow control mechanism allowing the sender to advance the transmission window linearly by one segment upon reception of an ACK, which indicates the last in-order packet received successfully by the receiver. When packet loss occurs at a congested link due to buffer overflow at the intermediate router, either the sender receives duplicate ACKs (DUPACKs), or the sender's retransmission timeout (RTO) timer expires. These events activate TCP's fast retransmit and recovery, by which the sender reduces the size of its congestion window to half and linearly increases as in congestion avoidance, resulting in a lower transmission rate to relieve the link congestion [11].

**TCP New Reno** is a prevalent TCP variant in the Internet. Its main components are the slow start algorithm and congestion avoidance phase. Slow start allows the sender to grow its window rapidly to capture bandwidth. The congestion avoidance phase following the slow start grows the congestion window linearly, promoting fairness and stability. Loss is detected either by a triple duplicate acknowledgment (ACK), in which case the congestion window is halved, or a retransmission timeout (RTO), in which case it is reduced to one segment. The congestion window evolves in Additive Increase Multiplicative Decrease (AIMD) fashion, which is appropriate for congestion-induced loss, but not for random loss due to wireless transmission errors [12].

**TCP Vegas** is a modified New Reno that detects loss proactively earlier; it can retransmit before receiving the third duplicate ACK or RTO timer expires. The congestion avoidance mechanism keeps the appropriate amount of data in the network, and congestion window increases and decreases in a linear fashion. Congestion window stability and low retransmission rate make Vegas a *good choice* for wireless networks.

**TCP Veno** specifically targets wireless networks, by combining Reno and Vegas mechanisms to deduce whether a network is congested or random

loss is more likely. Congestion window evolution is determined by throughput estimates, and reduction factors are 0.5 for congestion loss and 0.8 for random loss. Veno's efficient bandwidth utilization and low retransmission rate make it attractive for wireless links.

**TCP Westwood** enhances the window control and backoff process. Westwood relies on end-to-end rate estimation. The innovative idea is to continuously measure at the TCP sender the packet rate of the connection by monitoring the rate of returning ACKs while trying to find the bandwidth estimate, which is defined as the share of bottleneck bandwidth available to the connection. The estimate is then used to compute congestion window *cwnd* and slow start threshold *ssthresh* after a congestion episode (i.e., after three duplicate ACKs or a timeout). Westwood is a sender side modification of the congestion window algorithm aiming to improve the performance of Reno in wired as well as wireless networks. However, the available bandwidth estimation algorithm is complex and may not be able to keep up with the rapid changes in a hybrid wireless network.

**Cubic** is a TCP variant for networks with high bandwidth-delay product (BDP). It is characterized by an aggressive (cubic) window growth function that is independent of RTT. Upon loss, it reduces the window by a factor of 0.8. Cubic is the default TCP variant in newer Linux kernels.

**C-TCP**. With the idea that pure loss-based or delay-based congestion control approaches that improve TCP throughput in high-speed networks may not work well, this algorithm is designed to combine the two approaches. C-TCP can rapidly increase sending rate when a network path is underutilized, but gracefully retreat in a busy network when bottleneck queues grow. C-TCP is the algorithm included in Windows Vista and Windows Server 2008. However, due to the loss based component, CUBIC and C-TCP are not designed for high-loss wireless paths.

TCP	control	mobility support	modification requirements
TCP-Vegas	pro-active	low	sender side
TCP-NewReno	reactive	low	sender side
TCP-Veno	pro-active	low	sender side
TCP-Westwood	pro-active	high	sender side
TCP-Jersey	pro-active	high	sender and router side

Table 2. Comparison of TCP enhanced schemes in heterogeneous wireless communications.



The standard Reno scheme employs **reactive flow control**. The congestion window is adjusted based on the collective feedback of ACKs and DUPACKs generated at the receiver. TCP probes for the available bandwidth by continuously increasing the congestion window gradually until the network reaches the congestion state. In this sense, congestion is inevitable. TCP will then fall back to a much slower transmission rate, which may be unnecessary for wireless random loss [10].

In **proactive congestion control**, the sender attempts to adjust the congestion window proactively to an optimal rate according to the information collected via feedback, which can be translated to an indication of the network condition. By doing so, the sender reacts intelligently to the network condition or the cause of the packet loss, due to either congestion or random errors, and therefore prevents the network from entering an undesired state (e.g., congestion or unnecessary decrease of the congestion window). Different strategies can be employed in the design to provide the sender with the explicit network condition.

## 5.2 TCP traffic control

Classical methods for TCP traffic control are **priority queues, traffic shaping**, and instrumenting TCP packets to manipulate a receiver's advertised window: window resizing and Data/ACK Pacing [4].

- **Window Resizing.** TCP bandwidth estimation operates from the sender and cannot be easily manipulated. The advertised window of the receiver, however, can be decreased to reflect the actual bandwidth available in the wireless network. Limiting TCP sending behavior has beneficial effects even in the case when only TCP traffic is in the network. In order to control TCP sending rate without modification of TCP endpoints and maintain end-to-end semantics, the advertisement window could be modify in each ACK packet at the gateway. This method limits the total number of TCP data packets in transit between the endpoints. If the gateway changes the advertisement window based on the network status, TCP throughput can be limited close to its entry point. By keeping the window size small, retransmission and fairness problems among TCP flows are also relieved.
- **TCP Data/ACK Pacing.** Packet bursts is the problem that is not solved by window resizing. TCP pacing promises to reduce the burstiness of TCP traffic, and alleviate the impact of packet

loss, network delay, and delay jitter of VoIP traffic. TCP is pacing events out the transmission of a window of packets based on a shaper parameter  $R$ . After a packet of size  $pkt\_size$  goes out over the air, the next packet is scheduled no earlier than  $pkt\_size/R$ . The gateway chooses a rate  $R$  based on the network status to determine how much to send as well as when to send. One way to understand the impact of pacing is to consider burstiness from the perspective of network delay, jitter, and packet loss. With bursty traffic, packets arrive all at once at the gateway. As a result, queuing delay and delay jitter of VoIP packets grows linearly with TCP load due to large packet size, even when the load is below capacity. When carrying VoIP traffic at high load, 802.11 links are still perceived by TCP as being relatively free. This ignores the interference side effect that large TCP packets produce several hops away.

## 5.3 Available bandwidth estimation

To support VoIP in wireless networks, the classical TCP traffic control should be used in conjunction with methods to dynamically estimate available bandwidth in real time.

**TCP Reno** induces packet losses to estimate the available bandwidth in the network. While there are no packet losses, TCP Reno continues to increase its window size by one during each round trip time. When it experiences a packet loss, it reduces its window size to one half of the current window size. This is called *additive increase and multiplicative decrease*. The congestion avoidance mechanism adopted by TCP Reno causes a periodic oscillation in the window size due to the constant update of the window size. This oscillation in the window size leads to an oscillation in the round trip delay of the packets. This oscillation results in larger delay jitter and an inefficient use of the available bandwidth due to many retransmissions of the same packets after packet drops occur.

**TCP Vegas** adopts a more sophisticated bandwidth estimation scheme. It uses the difference between expected and actual flows rates to estimate the available bandwidth in the network. The idea is that when the network is not congested, the actual flow rate will be close to the expected flow rate. Otherwise, the actual flow rate will be smaller than the expected flow rate. TCP Vegas, using this difference in flow rates, estimates the congestion level in the network and updates the window size accordingly.

The **Westwood+** algorithm is based on end-to-end estimation of the bandwidth available along the TCP connection path. The estimate is obtained by filtering the stream of returning ACK packets and it is used to adaptively set the control windows when network congestion is experienced. In particular, when three DUPACKs are received, both the congestion window (*cwnd*) and the slow start threshold (*ssthresh*) are set equal to the estimated bandwidth (*BWE*) times the minimum measured round trip time (*RTT<sub>min</sub>*). When a coarse timeout expires the *ssthresh* is set as before, while the *cwnd* is set equal to one [12].

## 6 Conclusions

The major performance objectives of the WiMAX wireless communication systems are high link capacity and increase VoIP service from the user perspective. However, VoIP and TCP cannot simply share the WiMAX medium without severe voice quality degradation and/or reduction in TCP capacity. To evaluate the mobile WiMAX system capacity and performance, all the aspects of performance evaluation are presented: flat network architecture, prediction of system capacity, estimation of maximum number of VoIP users and throughput of TCP connection, VoIP and TCP traffic interaction, TCP enhancements and control tools.

The first WiMAX performance objective demands from TCP protocol to be aggressive in discovering wireless link bandwidth in order to achieve high utilization. However, there are no explicit signaling mechanisms in the network to tell TCP peers how fast to send, how much to send, or when to slow down a transmission. We presented research focused on the optimization of TCP performance in wireless networks. TCP enhanced variant Vegas uses the difference between expected and actual flows rates to estimate the available bandwidth in the network. Congestion window stability and low retransmission rate make *TCP Vegas* a good choice for wireless networks.

However, it is likely situation in WiMAX networks that none of the TCP endpoints can be controlled because upgrading TCP is infeasible or undesirable for other reasons. Even enhanced TCP endpoints cannot possibly protect wireless hops in the middle. We therefore presented the solution of control TCP traffic *before entering* WiMAX network in conjunction with methods to dynamically estimate available bandwidth in real time. Presented TCP control mechanisms could be used to limit the wireless resources taken by TCP but have different tradeoffs with respect to utilization and scalability.

## References:

- [1] K.R.Rao, Z.S.Bojkovic, D.A.Milovanovic, *Wireless multimedia communications: convergence, DSP, QoS and security*, CRC Press, 2009.
- [2] D.A.Milovanovic, Z.S.Bojkovic, Trends in multimedia over wireless broadband networks, *WSEAS Transactions on Communication*, Issue11, Vol.4, pp.1292-1297, 2005.
- [3] Z.S.Bojkovic, B.Bakmaz, Quality of Service and security as frameworks toward Next-Generation wireless networks, *WSEAS Transactions on Communication*, Issue 4, pp.147-152, 2005.
- [4] K.Kim, D.Niculescu, S.Hong, Coexistence of VoIP and TCP in wireless multihop networks, *IEEE Communication Magazine*, Vol.6, pp.75-81, 2009.
- [5] IEEE802.16 IEEE Standard for local and metropolitan area networks, Part 16: *Air interface for fixed broadband wireless access systems*, Amendment 2: *Physical and medium access control layers for combined fixed and mobile operation in licensed bands*, 2005.
- [6] G.Agapiou, et al., Experimental performance evaluation of the emerging WiMAX technology. *WSEAS Transactions on Information science and applications*, Issue 3, pp.322-327, 2006.
- [7] G.Leonardi et al, IEEE802.16e best effort performance investigation, in Proc. *IEEE International Symposium on Personal, Indoor and Mobile Radio Communications*, 2007.
- [8] Z.S.Bojković, D.A.Milovanovic, H.264 video transmission over IEEE802.11 based wireless networks: QoS cross-layer optimization, *WSEAS Transactions on Communication*, Issue 9, pp.1777-1782, 2006.
- [9] B-Ho Kim, et al., Capacity estimation and TCP performance enhancement over mobile WiMAX networks, *IEEE Communication Magazine*, Vol.6, pp.132-141, 2009.
- [10] Y.Tian, K.Xu, N.Ansari, TCP in wireless environments: problems and solutions, *IEEE Communication Magazine*, pp.27-32, 2005.
- [11] Yuan-Cheng Lai, Chang-Li Yao, Performance comparison between TCP Reno and TCP Vegas, *Computer Communications*, Vol.25, Issue 18, pp.1765-1773, 2002.
- [12] L.A.Grieco, S.Mascolo, Performance evaluation and comparison of Westwood+, New Reno, and Vegas TCP congestion control, *ACM Computer Communication Review*, 34(2), pp.25-38, 2004.