

Channel-Aware Earliest Deadline Due Fair Scheduling for Wireless Multimedia Networks

KHALED M.F. ELSAYED and AHMED K.F. KHATTAB

Department of Electronics and Communications Engineering, Cairo University, Giza, Egypt 12613
E-mails: khaled@ieee.org, ahmedkhatab@eng.cu.edu.eg

Abstract. Providing delay guarantees to time-sensitive traffic in wireless multimedia networks is a challenging issue. This is due to the time-varying link capacities and the variety of real-time applications expected to be handled by such networks. We propose and evaluate the performance of a channel-aware scheduling discipline and a set of policies that are capable of providing such delay guarantees in TDM-based wireless networks. First, we introduce the Channel-Dependent Earliest-Due-Date (CD-EDD) discipline. In this discipline, the expiration time of the head of line packets of users' queues is taken into consideration in conjunction with the current channel states of users in the scheduling decision. This scheme attempts to guarantee the targeted delay bounds in addition to exploiting multiuser diversity to make best utilization of the variable capacity of the channel. We also propose the violation-fair policy that can be integrated with the CD-EDD discipline and two other well-known scheduling disciplines [1, 2]. In this policy, we attempt to ensure that the number of packets dropped due to deadline violation is fairly distributed among the users. The proposed schemes can provide statistical guarantees on delays, achieve high throughput, and exhibit good fairness performance with respect to throughput and deadline violations. We provide extensive simulation results to study the performance the proposed schemes and compare them with two of the best known scheduling disciplines [1, 2] in the literature.

Keywords: fairness, multiuser diversity, QoS provisioning, scheduling, wireless networks

1. Introduction

The rapid growth of wireless technology, when coupled with the explosive growth of the Internet, has increased the demand for wireless data services. Traffic on beyond 3G wireless networks is expected to be a mix of real-time traffic such as voice, multimedia teleconferencing, and games, and data-traffic such as WWW browsing, messaging and file transfers. All of these applications will require widely varying and diverse quality of service (QoS) guarantees for the different types of offered traffic. Various scheduling disciplines have been developed in order to guarantee certain required QoS over wireline networks. However, these service disciplines, such as Weighted Fair Queuing (WFQ), virtual clock, Start-Time Fair Queueing (STFQ), and Earliest-Due-Date First (EDD) [3], are not directly applicable in wireless networks because they do not consider the characteristics of the wireless channel. These characteristics include high error rate, bursty errors, location-dependent and time-varying wireless link capacity, scarce bandwidth, user mobility, and power limitation of the mobile hosts.

All of the above characteristics make developing efficient and effective scheduling algorithms for wireless networks very challenging. Recent survey of the area of wireless scheduling

can be found in [4, 5]. Recently there has been increased interest in protocols for wireless networks which rely on significant interactions between various layers of the network stack. Generically termed cross layer design, many of these proposals are aimed at achieving performance improvements. For example, if wireless scheduling is performed based on the physical layer information (users' channel states), the efficiency of the wireless system in utilizing the system resources increases. This idea was firstly exploited by Knopp and Humblet [6] when they introduced a new diversity scheme, termed multiuser diversity. This type of diversity is inherent in a wireless network with multiple users sharing a time-varying channel. Multiuser diversity comes from the fact that different users usually have independent channel gains for the same shared medium (e.g. downlink). With multiuser diversity, the strategy of maximizing the total Shannon (ergodic) capacity turns out to be a greedy scheduling discipline where the scheduler allows at any time slot only the user with the best channel to transmit [7]. Results in [8] have shown that such a scheduling technique can increase the total (ergodic) capacity dramatically, in the absence of delay constraints, as compared to the traditionally used scheduling techniques.

One problem with such greedy scheduling is the unfairness in resource sharing among users in the network. This is due the fact that the users with the best channel conditions will always receive the biggest share of network resources; while the users suffering from bad channel conditions may not receive a reasonable share of the resources. The research reported in [9–14] was concerned with the problem of achieving throughput and/or temporal fairness among users. They report on scheduling schemes that attempt to guarantee that the difference in the services obtained by users with different channel conditions are as close as possible either on short-term basis and/or on long-term basis. However, these schemes provide no delay guarantees and thus are not suitable for delay-sensitive applications, such as voice and video.

In this paper we propose a scheduling discipline based on Earliest-Due-Date First (EDD) that exploits multi-user diversity and can provide statistical guarantees on delays, achieves high throughput, and exhibits good fairness performance with respect to throughput and deadline violations. We also present the violation-fair policy that can be integrated with the CD-EDD discipline and two other well-known scheduling disciplines [1, 2]. This policy ensures that the number of packets dropped due to deadline violation is fairly distributed among the users.

The rest of this paper is organized as follows: in Section 2 we describe the model of the system under consideration. Then in Section 3, a survey of scheduling delay-sensitive traffic in wireless network is presented. In Section 4, we propose the CD-EDD scheduling discipline and the violation-fair policy and describe their operation principles. In Section 5, a number of simulation experiments are carried out to investigate the performance of these schemes. Section VI summarizes the main findings of this paper.

2. System Model

We first describe the cellular wireless network model used, and more specifically the downlink of such a network. A base station transmits data to N mobile users, each of which requires certain QoS guarantees. In cell-structured wireless networks, the service area is divided into cells, and each is served via a base station. A single cell is considered in which a centralized scheduler at the base station controls the downlink scheduling, whereas uplink scheduling

uses an additional mechanism such as polling to collect transmission requests from mobile terminals [9, 10]. We assume that downlink and uplink transmission don't interfere with each other.

We consider a time slotted system, where time is the resource to be shared among users. A time-slotted cellular system can have more than one channel (frequency band), but at any given time, only one user can occupy a given channel within a cell. Here, we focus on the scheduling problem for a single channel over which a number of users could be time-division multiplexed. Time division multiple access (TDMA) systems divides the time into slots of length T_s , during which data transmission of a single user, the scheduled user, is carried out using all resources available to the base station at that time instant. For downlink scheduling, packets destined to different users are put in separate queues (in essence each flow is assigned a separate queue).

The time varying channel conditions of wireless links are related to three basic phenomena: fast fading on the order of milliseconds, shadow fading on the order of tens to hundreds of milliseconds, and finally, long time-scale variations due to user mobility. The channel fading processes of the users are assumed to be stationary, ergodic and independent of each other, and we also assume that the channel gains are constant over the slot duration. Since our algorithm will exploit the users' channel conditions in making the scheduling decision, we consider wireless systems with mechanisms to make predicted channel conditions available to the base station as is commonly the case with technologies such as HDR [15], UMTS-HSDPA [16], (E)GPRS [17], etc. The particular mechanism employed by a system depends on the communication standard. For example, in HDR and UMTS-HSDPA, the underlying physical channel uses explicit channel notification so that the scheduler has the best possible knowledge about the channel conditions. In UMTS-DCH, there is a logical control channel assigned with every user that allows a coarse estimation of the channel condition. The packet extensions (E)GPRS to GSM-TDMA systems offer various coding schemes to support data transmission over a wide range of channel conditions. These are typically switched on a slower timescale, e.g., based on experienced frame error rates. Regardless, the recently selected coding scheme can serve as a coarse indicator of the channel condition for the scheduler. In general, the faster and more precisely the channel quality can be predicted, the better the scheduler can incorporate this information into its decision as to which user to schedule next [9]. Thus, we will assume that the base station has the current (or delayed) channel state information of each user.

Figure 1 shows the architecture for channel-state aware scheduling of multiuser traffic over a fading time slotted wireless channel. The scheduler makes a decision to serve a particular queue at the beginning of every time slot. This decision could depend on HoL packet delay information, such as its waiting time and its time to expire, as well as channel states. Once a decision has been made, the chosen queue is serviced in that slot at the maximum possible rate corresponding to the state of its channel.

A very important and challenging problem in the design of high-speed communications networks is that of providing Quality of Service (QoS) guarantees, usually specified in terms of rate guarantees, loss probabilities or delays of packets in the network. The control of delays is often of crucial importance, especially for real-time applications such as audio and video streaming. Real-time traffic classes are modeled as a stream of packets, with each packet having an expiry time beyond which the packet is of no use to the end user. Expiry occurs when a packet has been waiting in the base station queue for a time greater than its deadline without being served. Such a packet is dropped by the system. The objective of the scheduler

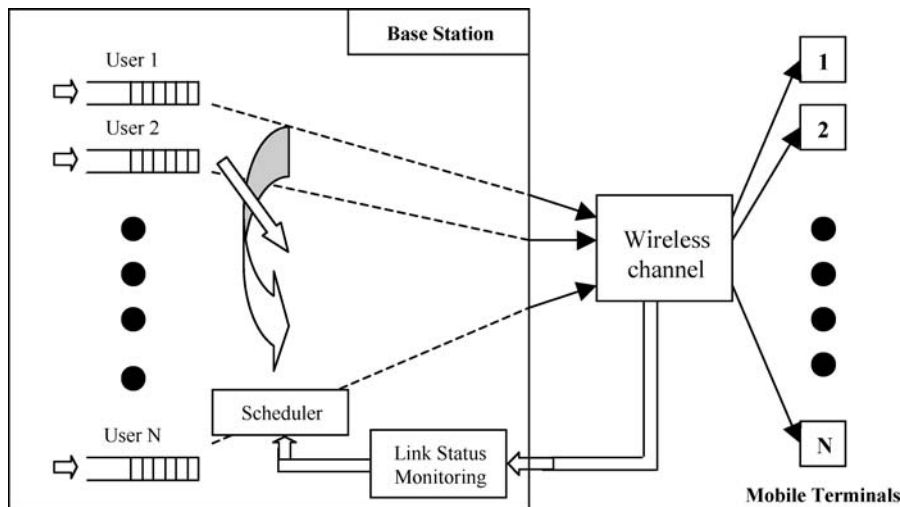


Figure 1. Channel state dependent downlink scheduling architecture for multiple users sharing a wireless TDMA channel.

is to transmit each packet before its expiry, and if this is not possible, to minimize the number of lost packets due to deadline expiry.

Such QoS requirements can be specified in terms of deadline T_i (deterministic QoS requirements) or accompanied with allowed violation probability δ_i (statistical QoS) for each user traffic flow. In this paper, we will use the following model in defining the QoS requirement of user i

$$P(W_i > T_i) \leq \delta_i \quad (1)$$

for $i = 1, 2, \dots, N$, where W_i is the delay encountered by user i packets. Under this constraint, the problem is to minimize the violation occurrences.

3. Previous Work

In this section, we present a survey of existing scheduling disciplines applicable in wireless networks with delay-sensitive users' traffic. General overview of the area of wireless scheduling can be found in [4, 5]. One fundamental issue in the wireless scheduling is the underlying channel model. For many proposals, the assumed channel model is a simple Markovian model where the channel can either be in a on (good) or off (bad) state. A continuous channel model is more realistic where according to the received signal-to-noise ratio and fading levels, certain rate can be achieved via adaptive coding and modulation at the physical layer. We start by providing a survey for schemes that are proposed for the on/off model and then proceed to those for the continuous channel model. Our work mainly assumes the continuous channel model, therefore no comparison with schemes for the on/off channel model is provided in the work.

3.1. WIRELESS SCHEDULING DISCIPLINES FOR THE ON/OFF CHANNEL MODEL

The Channel State Dependent Packet Scheduling (CSDPS) [18] is a wireless scheduling framework that allows the use of different disciplines such as round-robin, longest queue first, or earliest timestamp first. CSDPS does not have good fairness properties or provisions for throughput guarantees.

Idealized Wireless Fair Queueing (IWFQ) [19] is a realization of WFQ with a compensation mechanism for error-prone flows. Packets are tagged as in WFQ and is assigned a lead and lag counter indicating the state of the flow as compared to the error-free service. Packets are served exactly as in WFQ, however, when an error occurs, a flow with a bad channel is stopped and the scheduler chooses packets with the next smaller finish times from flows in good channel state. A flow losing service will be compensated after its channel becomes good. Compensation is guaranteed since flows returning from error states will have packets with smaller finish times. IWFQ enforces bounds on the leading and lagging counters such that the amount of compensation for lagging flows and penalty for leading flows are bounded. Throughput and delay guarantees can be provided. However, in IWFQ the main difficulty is in adjusting the values of the bounds since it exhibits a tradeoff between fairness and delay/throughput guarantees.

Channel-condition Independent Fair Queueing (CIF-Q) [20] is based on Start-Time Fair Queueing (STFQ). As in IWFQ, each flow is assigned a lead and lag counter that indicates the number of packets by which the flow is leading or lagging its error-free service. A leading flow relinquishes a parameterized portion of its assigned bandwidth to lagging flows. The CIF-Q discipline exhibits excellent performance with respect to delay bounds, throughput guarantees, and fairness. Both IWFQ and CIF-Q suffer from high algorithmic complexity since they need to simulate an error-free service and keep record of leading and lagging counters.

The Server-Based Fairness Approach (SBFA) is a wireless scheduling framework that can work with any underlying scheduling [21]. A portion of the bandwidth is specifically reserved for holding packets from flows suffering from bad-state and is termed the long-term fairness server (LTFS). When a flow is selected for transmission, it is allowed to transmit only if it experiences a good channel, otherwise a slot is inserted in the LTFS. When the HOL slot from the LTFS flow is chosen for transmission, the slot's original flow is allowed to transmit if it has a good channel, otherwise the next slot is selected. The problem with SBFA is that a flow with a good channel state may receive much more service than its promised share. Also, the LTFS must be pre-allocated resources for compensation.

Feasible Earliest Due Date (FEDD) [22] schedules the packet which has the earliest time to expire from the set of queues whose channels are marked in good state only. This algorithm showed unfairness in throughput sharing, since the user suffering from long periods of bad state will not be compensated by any mechanism, thus the authors suggested to use a rate-proportional scheduler which provides guaranteed minimum bandwidth to each connection and distributes the residual bandwidth using the FEDD criterion.

3.2. WIRELESS SCHEDULING DISCIPLINES FOR THE CONTINUOUS CHANNEL MODEL

The authors in [1] proposed a modification to the Largest Weighted Delay First (LWDF) [23] scheduling discipline that takes the time varying characteristics of wireless channels into account. The LWDF discipline is a parameterized version of first-input first-output (FIFO)

that works as follows: at the beginning of the time slot starting at time t , serve at the maximal possible rate the queue of user j , where

$$j = \arg \max_i \{a_i W_i(t)\} \quad (2)$$

where $W_i(t)$ is the waiting time of head of line (HoL) packet of the i th user at the time slot starting at time t , and $a_i > 0$, $i = 1, 2, \dots, N$, are constant weights. If the delay QoS requirements for all users is as expressed by (1), it was proved in [23] that the choice of weights a_i that makes LWDF discipline nearly throughput optimal¹ (note however that this choice of weights is valid only for large values of the delay bound T_i and very small values of δ_i) is:

$$a_i = \frac{-\log(\delta_i)}{T_i} \quad (3)$$

The proposed modification of [1] was to exploit multiuser diversity in order to increase the efficiency of channel utilization (and hence the system throughput) and also compensate delayed users. The proposed Modified Largest Weighted Delay First (M-LWDF) discipline schedules the j th user, where

$$j = \arg \max_i \{\gamma_i \mu_i(t) W_i(t)\} \quad (4)$$

where $\mu_i(t)$ is the state of the channel of user i at time t , i.e. the actual rate supported by the channel. This rate is assumed to be constant over one slot. It has been proven in [23] that setting $\gamma_i = a_i / \bar{\mu}_i$, where a_i is given by (3) and $\bar{\mu}_i$ is the mean rate supported over the i th channel, makes the M-LWDF discipline throughput optimal. In practice, the mean rate can be measured over a certain, but relatively long, time window [8] by averaging the rate actually given to that user in that window. The M-LWDF scheduling discipline could be rewritten as schedule the j th user, where

$$j = \arg \max_i \left\{ a_i \frac{\mu_i(t)}{\bar{\mu}_i} W_i(t) \right\} \quad (5)$$

The discipline provides good QoS for delay sensitive users only with properly chosen parameters and can be easily implemented. It was shown in [1] that M-LWDF discipline is throughput optimal. It was also shown how M-LWDF can be used to achieve alternative QoS defined in terms of a predefined minimum long-term throughput for each user. Unfortunately, M-LWDF scheduling was found to be highly dependent on the value of the parameters a_i , and its performance changes significantly with the QoS requirements of the users' flows.

In order to reduce the dependency of the M-LWDF discipline on the settings of the parameters a_i , the authors in [1] also proposed a new scheduling discipline, which was further investigated and implemented in [2] for CDMA/HDR system and further modified in [24]. The proposed scheme was called the exponential rule scheduling discipline. This discipline

¹ A discipline is called throughput optimal if it can handle all the offered traffic and render the stability of all queues if this is feasible to any other discipline, i.e. has the largest stable admission region.

schedules the j th user at the time slot starting at time t for transmission, where

$$j = \arg \max_i \left\{ a_i \frac{\mu_i(t)}{\bar{\mu}_i} e^{\frac{a_i W_i(t) - \overline{aW}}{1 + \sqrt{\overline{aW}}}} \right\} \quad (6)$$

and

$$\overline{aW} = \frac{1}{N} \sum_{i=1}^N a_i W_i(t) \quad (7)$$

This discipline attempts to equalize the weighted delays $a_i W_i(t)$ of all the queues when their differences are large. If one of the queues would have larger (weighted) delay than the others by more than order $\sqrt{\overline{aW}}$, then the exponent term becomes very large and overrides the channel considerations (as long as its channel can support a non-zero rate), hence giving priority to that queue. (It can be easily noticed that the \overline{aW} term in the exponent can be dropped without changing the rule since it is common for all queues. This term is present only to emphasize the motivation of the rule). On the other hand, for small weighted delay differences (i.e. less than order $\sqrt{\overline{aW}}$), the exponential term is close to unity, and the discipline behaves as the proportionally fair discipline. Hence, the exponential rule discipline gracefully adapts from a proportionally fair discipline to a one that balances delays [1, 2]. It was proven in [25] that the exponential discipline is throughput optimal with general assumptions on the channel and arrival processes. Moreover, simulation results in [2] showed that the exponential rule scheduling discipline exhibits better delay tails compared to any other scheduling discipline in the sense that the delays of all users are about the same and are all reasonably small. This occurs, however, for large values of T_i and very small values of δ_i , which is not desired practically. Moreover, as in the M-LWDF discipline, the exponential rule scheduling discipline was found to be highly dependent on the parameter settings. Therefore, it was advised in [1] that identifying good scheduling disciplines which are less dependent on the “proper” parameter setting would be desirable.

4. The Proposed Scheduling Schemes

The goal of our work is to design scheduling schemes for delay sensitive traffic that exhibit “good performance” and that exploits multi-user diversity inherent in wireless communications. Based on the reported shortcomings of previous work reported in the literature as outlined in Section 3, we mean by the word “good performance” that such a discipline should attempt to achieve the following objectives:

1. *Maximize the overall system throughput:* This could be easily achieved if the scheduling discipline utilizes multiuser diversity in order to efficiently utilize the channel capacity by giving higher priority to the user with the best channel conditions at a certain time instant, which means that this user can transmit with the highest possible rate, and thus increase the system throughput.
2. *Graceful compensation of large delays:* For a real-time traffic packet, it is necessary that its delivery be done within its deadline, otherwise, the information contained in this packet will be irrelevant for the receiver. Thus, a good scheduling discipline should have a mechanism

to compensate queues whose packets are experiencing long delays in the system in order to prevent their packets from being dropped. This case may be encountered by users far from the base station resulting in channels suffering from long periods of bad conditions. Such a mechanism will guarantee the QoS requirements if defined as delay bound or packet loss ratio. Thus it minimizes the number of packets dropped due to deadline violation, which in turn increases the system throughput.

3. *Fairness in resource sharing:* Fairness is an intuitively desirable property of scheduling disciplines. A fair scheduling discipline should distribute the resources available to the system, such as capacity and time, fairly among different users (of the same service class). Fairness may be accomplished in delay distributions, service rates, number of packets lost, etc.
4. *Weak dependency on the parameters setting:* As advised in [1], it is desirable to identify good scheduling disciplines which are less dependent on the proper parameter setting. In other words their performance does not change significantly when the QoS requirements are widely changed.

4.1. THE CHANNEL DEPENDENT EARLIEST DUE DATE (CD-EDD) SCHEDULING DISCIPLINE

In classical wireline earliest due date scheduling, each packet is assigned a deadline, and the scheduler serves packets in order of their deadlines. The queue with the smallest deadline is served first by the maximum available rate. If the scheduler is overcommitted, then some packets miss their deadlines. EDD cannot be efficiently employed in wireless networks since it does not consider the time varying characteristics of wireless links. It was not reported in the literature the existence of a scheduling discipline that combines the EDD scheduling concept with a mechanism to adapt with the characteristics of wireless networks (with the exception of the attempt in [22] which does not actually adapt with the time varying nature of wireless channel, since it assumed the simplified two-state Markovian channel model. Also, while M-LWDF takes delay into account, it is not based on the EDD mechanism).

We propose a new scheduling discipline, which we call the channel dependent earliest due date first (CD-EDD) discipline. This is basically a channel-state dependent EDD discipline where the scheduler chooses to schedule the queue whose HoL packet has the earliest time to expire and the best channel conditions, and consequently the highest transmission rate, among all queues. The proposed CD-EDD scheduling discipline is as follows:

At the time slot starting at time t , schedule with the maximum possible rate the queue of the j th user, where

$$j = \arg \max_i \left\{ a_i \frac{\mu_i(t) W_i(t)}{\bar{\mu}_i d_i(t)} \right\} \quad (8)$$

where a_i is the weighting parameter reflecting the statistical QoS requirements of the i th user. $\mu_i(t)$ is the actual rate that could be used for transmission by the i th user at time t , which reflects the current channel state of the user's channel. $\bar{\mu}_i$ is the mean rate supported or previously offered to the i th user. $W_i(t)$ is the delay experienced by the HoL packet since its entrance to the i th user queue in the base station. $d_i(t)$ is the time to expire of the i th user HoL packet, which is the difference between the deadline, T_i , and the delay experienced till time t , $W_i(t)$,

i.e.

$$d_i(t) = T_i - W_i(t) \quad (9)$$

The behavior of the CD-EDD discipline can be explained as follows: when a certain queue has its HoL packet waiting in the system for a relatively long period (but have not expired yet), its time to expire will decrease significantly. In such a situation, the term $W_i(t)/d_i(t)$ will grow significantly due to the contribution of $1/d_i(t)$ until it overcomes other terms in (8). This has an effect akin to reducing the number of dropped packets due to deadline violation. On the other hand, if the delay characteristics of all users are about the same, i.e. their time to expire and waiting times are close, the term $W_i(t)/d_i(t)$ will be common to all users, and the discipline then reduces to a proportionally fair scheduler that exploit multiuser diversity to efficiently utilize the channel bandwidth of multiuser systems in a fair manner. It is worth mentioning that weights a_i doesn't contribute significantly in the decision. A rule of thumb for choosing a_i which works in practice is the one given in (3) since this choice is based on large deviations optimality results.

The CD-EDD is a scheduling discipline that can be used to provide QoS guarantees, defined in terms of delay bounds, for real-time traffic in wireless networks. This is achieved by increasing the priority of delayed users to get fair access to the medium over time. An important feature of the CD-EDD discipline is its weak dependency on the value of deadline required, and thus can be used for a wide variety of QoS requirements.

4.2. THE VIOLATION-FAIR POLICY

Another new idea than can be applied in conjunction with any scheduling discipline in order to enhance their fairness characteristics is proposed here. This is based on the number of deadline violations occurring to packets of different queues. This requires that each queue in the base station be accompanied with a counter that counts the number of packets lost in this queue's flow. This may be implemented practically by means of sliding window basis. Let us define $NV_i(t)$ to be the number of deadline violations encountered in the flow of the i th user up to time t , and $\overline{NV}(t)$ to be the average of the number of violations in all N queues up to time t , i.e.

$$\overline{NV}(t) = \frac{1}{N} \sum_{i=1}^N NV_i(t) \quad (10)$$

The scheduler may use the number of deadline violations to find a way to compensate users suffering from unfairness in the number of dropped packets. For example the scheduler could give more credit or increase the priority level so that such a user could access the system resources. This could be achieved by a scheduling discipline that utilizes a term like $NV_i(t)/\overline{NV}(t)$ in making the scheduling decision. We call such a scheduling discipline a violations-fair (VF) discipline.

Initially, we applied this idea directly to the proportionally fair [8] scheduling discipline, yielding the violations-fair proportionally fair discipline. So, in each time slot the scheduler

chooses the j th user, where

$$j = \arg \max_i \left\{ \frac{\mu_i(t) NV_i(t)}{\bar{\mu}_i \overline{NV}(t)} \right\} \quad (11)$$

This reduces the number of packets lost for users with bad channel conditions, therefore, enhancing the fairness characteristics of the proportionally fair discipline. On the other hand, it still lacks a mechanism for provisioning of QoS guarantees for delay sensitive traffic.

We further apply the proposed modification to both the M-LWDF and the exponential rule disciplines. In the violations-fair modified largest weighted delay first (VF-M-LWDF) discipline, the scheduling decision is such that at the time slot starting at time t , schedule with the maximum possible rate the queue of the j th user, where

$$j = \arg \max_i \left\{ a_i \frac{\mu_i(t) NV_i(t)}{\bar{\mu}_i \overline{NV}(t)} W_i(t) \right\} \quad (12)$$

In the violations-fair exponential (VF-EXP) discipline, the scheduling decision is such that at the time slot starting at time t , schedule with the maximum possible rate the queue of the j th user, where

$$j = \arg \max_i \left\{ a_i \frac{\mu_i(t) NV_i(t)}{\bar{\mu}_i \overline{NV}(t)} e^{\frac{a_i W_i(t) - aW}{1 + \sqrt{aW}}} \right\} \quad (13)$$

The proposed modification enhance their performance since the addition of the violations-fair term will ensure fairness in both the delay times and throughput. This could be explained since both the M-LWDF and the exponential disciplines minimize the packet delay, and when the number of dropped packets is fairly distributed among the users, the long term service rate will be equal for all users. Another very important gain of the violations-fair version of these disciplines, is that their performance is not much dependent on the parameter setting as was the case in the original disciplines (this will be seen in the simulation results).

Finally, if the proposed violations-fair technique is applied on the proposed CD-EDD discipline, we can get a scheduling discipline applicable in wireless network that explicitly provides QoS to delay sensitive traffic, with excellent fairness characteristic with respect to data rate, delay bound, and delay bound violation. The violations-fair-channel-dependent earliest deadline due date (VF-CD-EDD) scheduler chooses, at the time slot starting at time t , the j th user, where

$$j = \arg \max_i \left\{ a_i \frac{\mu_i(t) W_i(t) NV_i(t)}{\bar{\mu}_i d_i(t) \overline{NV}(t)} \right\} \quad (14)$$

In the next section, we provide an extensive set of simulation results that explore the performance of the proposed CD-EDD discipline compared with the M-LWDF and the exponential rule disciplines as a reference. We also show the advantages achieved by their violations-fair versions, namely the VF-CD-EDD, the VF-M-LWDF, and the VF-exponential scheduling disciplines.

5. Performance Evaluation

5.1. SIMULATION SETUP

The system model used in simulating the wireless cell-structured channel-aware scheduler was described in Section 2. We chose the High Data Rate (HDR) CDMA system model. HDR technology has recently been proposed as a TDM-based overlay to CDMA with the goal of providing packet data services to mobile users. HDR is a downlink packet data service that occupies a single data carrier of a CDMA system, where users share the channel in a time division multiple access manner, i.e. at each time slot T_s only one user can transmit its data at the full power available to the base station. A very attractive feature of HDR is enabling the use of efficient scheduling algorithms since it provides a mechanism for link status monitoring as described in Section 2.

The cell serves N mobile users each receiving a data flow. The base station contains N queues, each corresponding to a different data flow and an associated scheduler. The scheduler makes a decision every 1.667 millisecond based on the current information available at the start of the time slot. As we are mainly interested in scheduling users with time sensitive traffic, we model the packet arrival processes to each of the N user's queues as a Bernoulli processes with a mean rate of 28.8 Kbps. This rate corresponds to the typical rate required for streaming audio over the Internet. Real time users, like streaming audio, will indeed generate a smooth traffic, and hence, a Bernoulli model seems reasonable for such traffic. Like the original EDD, the CD-EDD is expected to be throughput optimal for such traffic model. The HDR packet size is 128 bytes. The QoS requirements of each user are expressed in the form of the probability that the waiting time encountered by a typical packet of the i th user stream exceeds the deadline T_i is less than or equal to δ_i . We assume for simplicity that all users require the same service quality, i.e. they all have the same T_i and δ_i .

Even though all users share a common channel, the channel capacity, of that channel seen by different users is different. This is due to the wireless link characteristics described above. The instantaneous capacity of a wireless channel is given by

$$C(t) = B \log_2(1 + |h(t)|^2 SNR) \quad (15)$$

where $C(t)$ is the channel capacity or the data rate (in bits per second) that can be transmitted on a channel of bandwidth B (Hertz). The bandwidth of HDR/CDMA channel is 1.25 MHz. The term $|h(t)|$ is the normalized gain (or fading level) of the wireless channel at time t , and SNR is the required signal to noise ratio at the receiver antenna (13dB for HDR/CDMA system). For simulation purposes, we use the typical HDR and cell parameters given in [15]. The average fade level distribution of a typical mobile in HDR cell can be easily found. The fading process of each user's channel can be represented by a Rayleigh process. So, in order to simulate N channels, we pick N fading levels according to the above distribution, and generate N Rayleigh processes with means equal to these fading levels after being normalized. A sample of such a process represents the normalized channel gain of a certain channel at time t with a probability $p_R(r) = \frac{r}{\sigma^2} e^{-\frac{r^2}{2\sigma^2}}$, $r \geq 0$, where σ is related to the mean normalized fading level $|\bar{h}|$ according to, $|\bar{h}| = \sigma \sqrt{\frac{\pi}{2}}$. Accordingly, the mean data rate $\bar{\mu}_i$ that can be supported on the a

channel of mean fading level $|\bar{h}_i|$ is

$$\bar{\mu}_i = B \log_2 (1 + |\bar{h}_i|^2 SNR) \quad (16)$$

The user with the highest value $\bar{\mu}_i$ will have the best channel condition.

It is worth mentioning that HDR doesn't support continuous transmission rates, i.e. the scheduled user cannot transmit with the rate computed above but with the maximum possible rate from a set of discrete rates. An HDR user can transmit data at a rate of $9.6 * 2^i$ Kbps, $i = 0, 1, \dots$, with a maximum rate of 2 Mbps. Thus the state of channel $\mu_i(t)$ at the start of the time slot at time t will be the actual rate that the channel can support, rather than the channel capacity at that time instant. We have assumed that the channel conditions do not change significantly within a time slot duration. Finally, all simulations will be carried out for a duration of 10 minutes.

5.2. PERFORMANCE METRICS

Here, we discuss the performance metrics used to evaluate the performance of the various scheduling algorithms. These are the delay, throughput, and packet loss criteria. We briefly outline them in turn.

For real-time traffic, a good measure of performance is the delays packets incur at the base station. A good scheduling algorithm should keep all delays below the delay bound T_i with high probability. The delay distribution curves can be used to illustrate the delay behavior of the scheduling disciplines under consideration. As remarked earlier, scheduling algorithms which keep the delays of all the users about the same and keep them all reasonably small are superior to those which may have better delay tails for one of the users but have very bad delays for other users.

Some parameters of the delay distribution, such as the worst-case delay, the mean delay, or the 95-percentile delay, could be used to evaluate the QoS received by a user. We will consider the 95-percentile delay as our measure of the delay guarantees offered by a scheduling discipline. The 95-percentile delay is defined as the value of the delay where ninety-five percent of the users' packets experience delays smaller than that value.

Unlike non-real-time users, who may have their QoS requirements in the form of a guaranteed minimum rate, real-time users do not need the scheduler to preserve certain bandwidth for their packets' transmission. So that, we will only take the total throughput achieved by the system as a measure of the throughput performance of the scheduling disciplines at hand. The fairness of bandwidth sharing among users is also a good indication of the efficiency of any scheduling discipline.

The fraction of packets dropped, due to deadline violation, for a user can be used to evaluate the loss performance of a scheduling discipline. This fraction is required to be as small as possible in order to say that the scheduler is suitable for scheduling real-time traffic. From fairness point of view, it is better to equalize the fraction of packets lost in different queues.

5.3. RESULTS AND DISCUSSIONS

First, we estimate the number of users, with traffic like the one described above, that can be supported by a single HDR cell under each of the previously mentioned scheduling policies. In

Table 1. System capacity for different delay bounds

Delay bound (m sec)	20	60	100	200	300	400
CD-EDD	4	12	16	16	16	16
M-LWDF	4	10	12	16	16	16
EXP	4	8	12	16	16	16

this experiment, we simulate N users, uniformly distributed throughout the cell, and monitor the average service rate received by a single user for different values of N . As N increases, and as long as the channel capacity can support such a number of users, it is expected that the service rate received by each user will be close to its arrival rate. Any further increase of the number of users, while keeping the channel capacity unchanged, will make the scheduler unable to serve additional users in the appropriate time so more packets will be dropped and thus the average throughput share of each user will decrease. So, we will take the number of users beyond which the average service rate received by any user in the system start to decrease as the system capacity.

Table 1 summarizes the simulation results of system capacity achieved by different disciplines for different values of the delay bound and for a violation probability (δ_i) of 5% (which will be used for all experiments). It is clear that the CD-EDD discipline provides better performance than other disciplines.

In the second experiment, we investigate the delay performance of various scheduling disciplines. We begin with a comparison of the proposed CD-EDD scheduling disciplines versus both the M-LWDF and the exponential rule scheduling disciplines reported as the most suitable policies for scheduling delay sensitive traffic in the literature. In order to be able to evaluate the performance of the scheduling techniques, the system should be loaded with its maximum capacity. Based on the results of the first experiment in Table 1, we assume that the cell is serving 14 mobile terminals, i.e. $N = 14$, and generate 14 i.i.d. Bernoulli processes each with a mean rate of 28.8 Kbps. Using the procedure described before, we also generate 14 Rayleigh-faded channels for the users uniformly distributed in the cell.

In Figure 2, we plot the delay distribution tails for both user 1 (with the best channel conditions) and user 14 (with the worst channel conditions) for the M-LWDF, exponential rule, and the proposed CD-EDD scheduling disciplines for delay bounds of 60, 100, 200, and 400 milliseconds.

These bounds are encountered practically in multimedia streams. It is obvious that all policies will have the same performance for very small bounds, e.g. 20 milliseconds, with such a small deadline, the system cannot serve such a number of users with this high QoS. We observe that for moderate delay bounds, e.g. 60 and 100 milliseconds, the performance of the exponential rule scheduling slightly outperforms both the M-LWDF and CD-EDD disciplines. In this case, all delays are kept at small value and about the same for all users. For higher bounds, e.g. 200, and 400 milliseconds, we find that the delay performance of the CD-EDD scheduler does not change significantly. On the other hand, the performance of both the M-LWDF and the exponential rule schedulers degrades severely as the gap between the tails of the best user and the worst user becomes wider. This severe degrading in performance is caused by the dependency of both disciplines on the quality of service required, which affects the value of the weights $\{a_i\}$ that controls the performance of these disciplines. (Keep in mind that the goal of the M-LWDF is to minimize the weighted delays while the exponential rule

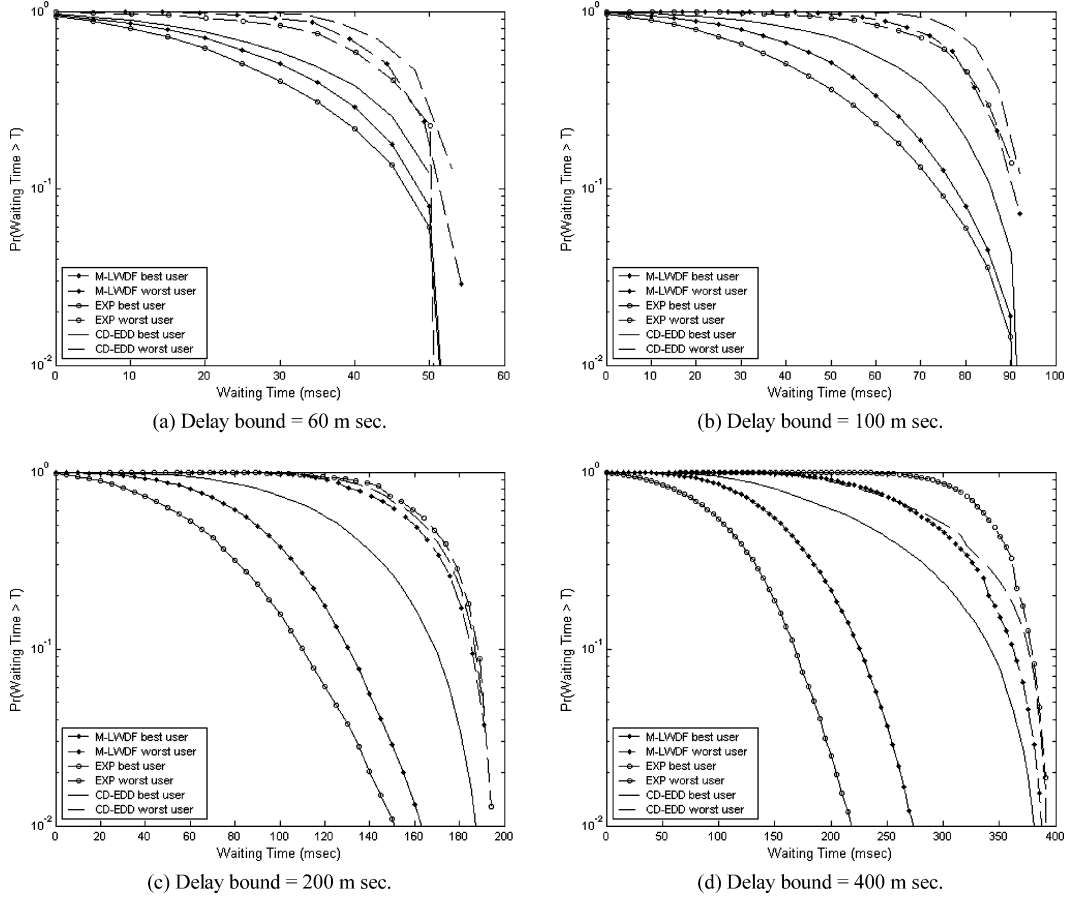


Figure 2. Delay distributions of the best user and worst user for different delay bounds.

scheduling tries to keep all the delays around the average of these weighted delays (\overline{aW}), and they both do not target a certain delay bound.) On the other hand, the CD-EDD discipline does not suffer from such a dependency on the QoS (and consequently the weight values), since the philosophy of this discipline is mainly to serve more packets before their deadlines expire. This may cause the packets of the best channel users to experience relatively higher delays than in the case of other disciplines, but still lower than the worst user.

This can be further demonstrated when we plot the maximum and the minimum (corresponding to the best and worst users) 95-percentile delay and percentage of packets dropped due to deadline violations versus the delay bounds as shown in Figure 3a and b, respectively. It is clear that it is not desirable to keep one or some users' delays below a value much smaller than the required bound while leaving one or some users suffering from dropping a large percentage of their packets as the case with both the M-LWDF and exponential rule schedulers. While in the CD-EDD discipline, the maximum and the minimum 95 percentile delays are about the same and so close to the delay bound, besides the packet loss ratios are very small and very close to each other. So, the base stations can guarantee strong delay bounds for all delay sensitive users in a fair manner by using the CD-EDD scheduling discipline, regardless of the value of these bounds.

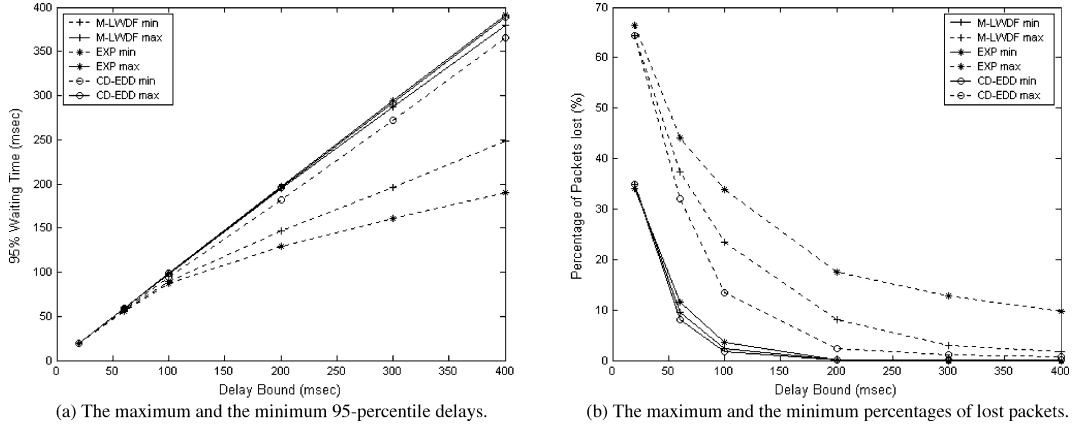


Figure 3. Fairness behavior of the original disciplines.

When the same experiments were carried out for the violations-fair versions of the above disciplines, it was found that the performance of all the violations-fair policies is not much dependent on the QoS required (including the disciplines which was originally suffering from that dependency). Figure 4 shows the delay distribution tails of the best user and the worst user for different delay bounds. It is observed that the tails became much more closer for the VF-CD-EDD, but with a little bit higher delays than those of either the VF-M-LWDF or the VF-EXP policies. As shown in Figures 5a and b, where the maximum and the minimum 95-percentile delay and percentage of packets dropped due to deadline violations are plotted versus the delay bounds, the service quality offered to different users, in terms of 95% delay, is about the same. Furthermore, the amount of packets dropped in the system due to deadline violations becomes very small. Moreover, this amount is distributed among all users in a fair manner. Like the CD-EDD, the VF-CD-EDD have the superiority since it achieves the smallest number of packets lost due to deadline violations.

Finally, we study the throughput characteristics of the aforementioned scheduling disciplines. Here we are interested in studying the overall throughput of the system as well as how the throughput is divided among users with different channel conditions, i.e. fairness in throughput sharing. The results of this experiment are illustrated in Figure 6 for 100 milliseconds bounds by plotting the overall throughput of the system versus the number of users using the simulation setup of the first experiment. When the system is operating with number of users less than the system capacity, the total throughput of the system equals the sum of the arrival rate. When the system is serving more users than the system capacity, we observe that, regardless of the delay bound, the total throughput achieved using the CD-EDD is higher than any other discipline because it causes the smallest number of packet to be lost due to deadline expiry. Moreover, the total throughput achieved introducing the violations-fair policy is less than the throughput achieved with the non-violations fair counterpart. This is because in order to achieve fairness in the ratio of packets dropped among different user, the violations-fair policies may prevent users with good channel conditions from transmitting their data for the sake of users with bad channel conditions (such channels support low transmission rates only). So the overall throughput achieved by the system will be lower than the case where the scheduling discipline do not intend to make users with good channel condition drop some packet for the purpose of fairness.

Table 2. Throughput fairness index for a 100 m sec delay bound.

	CD-EDD	EXP	M-LWDF
Original discipline	0.8749	0.6781	0.7797
Discipline with violations-fair policy	0.9656	0.9386	0.9520

In order to evaluate the throughput fairness performance of the proposed scheduling schemes, let us define the throughput fairness index as the ratio between the lowest achieved rate and the highest achieved rate. The more the fairness index approaches unity, the better the scheduling discipline in the sense that all users almost have been served with the similar rate as long as they have the same QoS. Table 2 lists the computed values of the throughput fairness index for both the original scheduling disciplines as well as their violations-fair counterparts. We use the same simulation model used in the previous experiment where the base station serves 14 users.

Concerning the non-violations fair disciplines, it was found that the CD-EDD discipline provides the highest fairness index compared to both the exponential rule and the M-LWDF

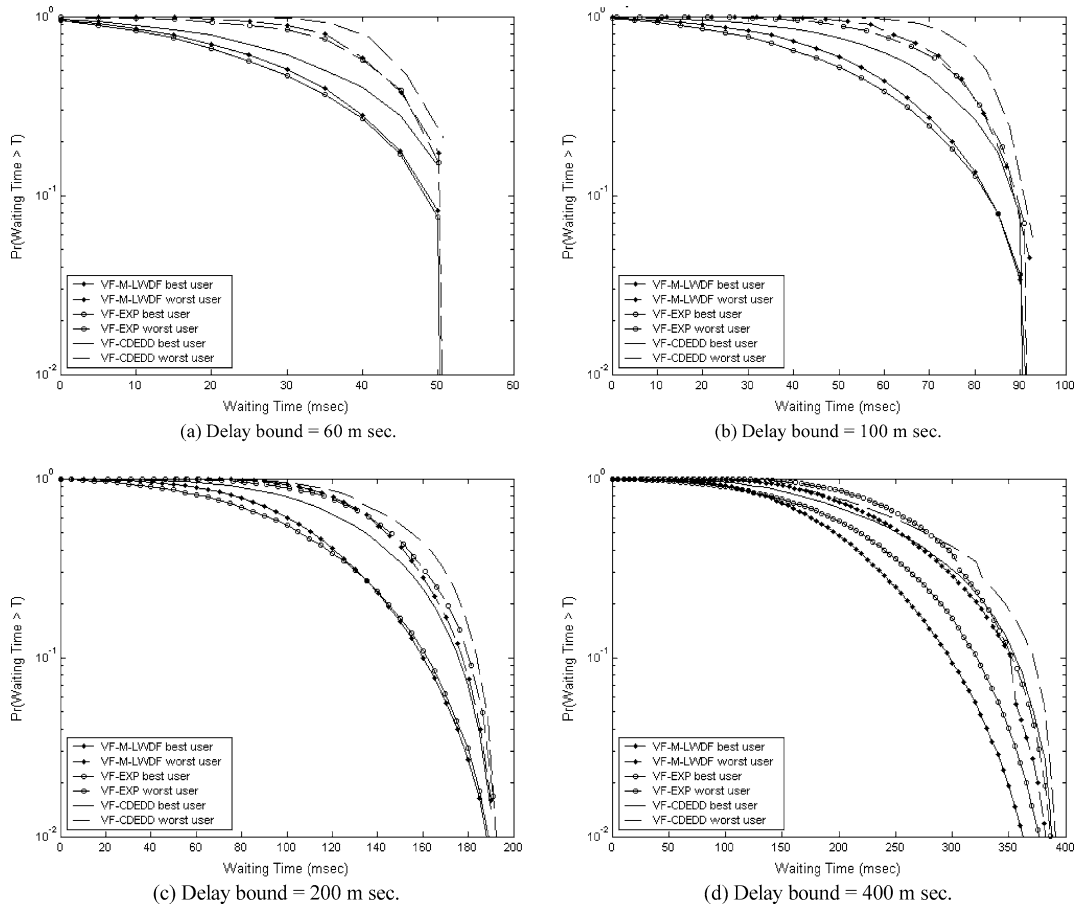


Figure 4. Delay distributions of the best user and worst user for different delay bounds for the violations-fair policies.

Channel-Aware Earliest Deadline Due Fair Scheduling

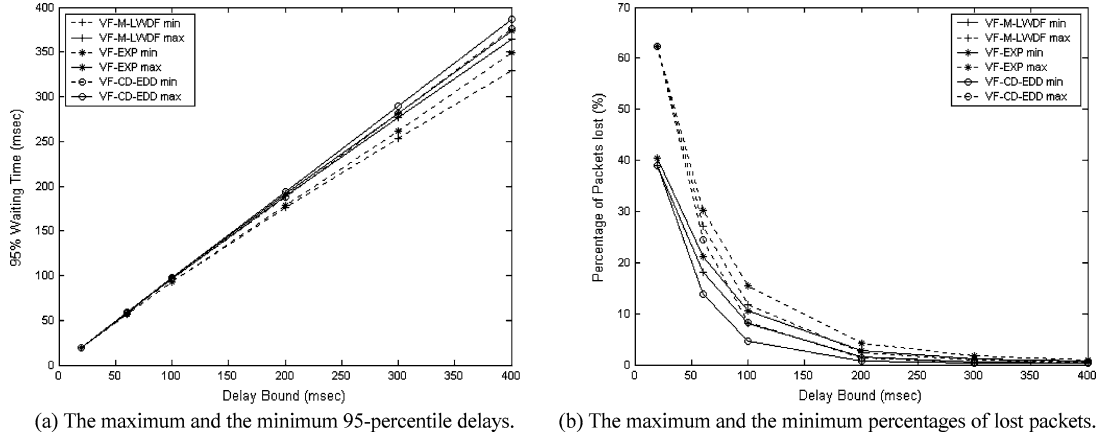


Figure 5. Fairness behavior of the violation fair disciplines.

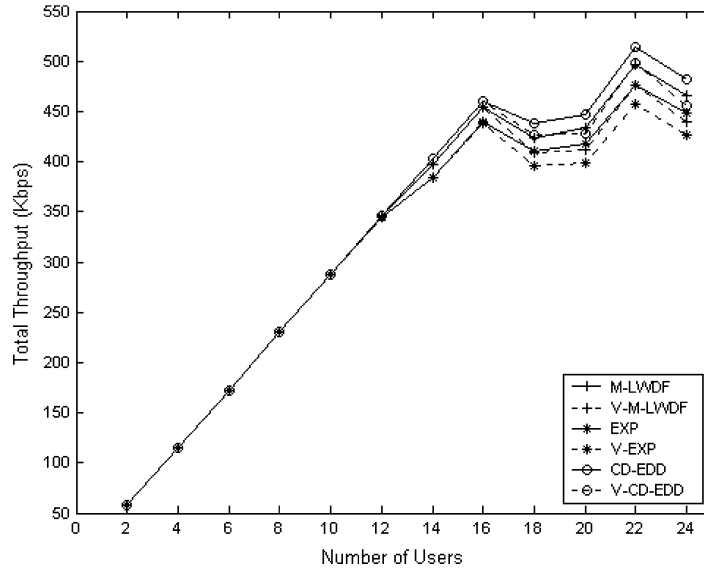


Figure 6. The total throughput achieved by different disciplines at 100 m sec delay bound.

disciplines regardless the delay bound (we only include the results at 100 m sec for illustration purposes). On the other hand, even though the violations-fair disciplines leads to a slightly lower throughput per user as previously discussed, they ensure fairness in throughput sharing among all users regardless of the delay bound for the three considered violations-fair disciplines. This is because in such disciplines, the number of packets dropped due to deadline expiry is almost equal for all users. However, we note that the VF-CD-EDD has slightly better throughput fairness performance. We summarize the main results of our simulation experiments in Table 3.

6. Conclusions

This paper addresses the problem of scheduling real-time users over TDM-based wireless multimedia networks. We introduced the Channel Dependent Earliest-Due-Date first (CD-EDD)

Table 3. Simulation results summary

	CD-EDD	M-LWDF	EXP	VF-CD-EDD	VF-M-LWDF	VF-EXP
Delay bound guarantee	✓	X	X	✓	✓	✓
Independency of QoS requirement	✓	X	X	✓	✓	✓
Total system throughput	Highest	Moderate	Lowest	Less than original counterparts		
Delay fairness	✓	Depend on QoS		✓	✓	✓
Throughput fairness	Best	Good	Worst	✓	✓	✓
Delay bound violation fairness	✓	Bad	Worst	✓	✓	✓

scheduling discipline, a discipline that attempts to provide statistical delay bound guarantees for time-sensitive traffic in networks with time-varying channels. Gains in throughput and realized delay are achieved by exploiting multi-user diversity techniques in which the scheduling decision takes into account the current channel state for each user in the system. By considering the packets dropping due to deadline violation, we also presented a set of scheduling policies that has satisfactory fairness characteristics in delays, throughput, and packet loss ratios among different users regardless of the value of the delay bound.

Simulation results of the proposed schemes showed that the services received by different real-time users, namely, delays, rates, and loss ratios, can be fairly achieved for a wide range of applications. The proposed schemes outperform other existing disciplines. The computational complexity of the proposed schemes are low and are suitable for application in future broadband fixed or mobile wireless systems such as 802.16a and 802.20.

It is known that using stringent deadline scheduling policies in a wireless system also increases the probability that packets are scheduled when the channels are in bad conditions. This could lead to some packet loss and retransmission delays. We believe that a good extension for future work would be studying the tradeoff between good application layer throughput and deadline violation and finding the best balance point. Also, whether the proposed CD-EDD discipline is throughput-optimal remains an open question.

References

1. M. Andrews, K. Kumaran, K. Ramanan, A.L. Stolyar, R. Vijayakumar, and P. Whiting, "CDMA Data QoS Scheduling on the Forward Link with Variable Channel Conditions", Bell Laboratories Technical Report, April, 2000.
2. S. Shakkottai and A.L. Stolyar, "Scheduling Algorithms for a Mixture of Real-Time and Non-Real-Time Data in HDR", in *Proceedings of the 17th International Teletraffic Congress - ITC-17*, Salvador da Bahia, Brazil, pp. 793–804, 24–28 September, 2001.
3. S. Kashave, *An Engineering Approach to Computer Networking: ATM Networks, the Internet, and the Telephone Network*, Addison Wesley, 1997, ch. 9.
4. Y. Cao and V. O. K. Li, "Scheduling Algorithms in Broadband Wireless Networks", *Proceedings of the IEEE*, pp. 76–87, Jan. 2001.
5. H. Fattah and C. Leung, "An Overview of Scheduling Algorithms in Wireless Multimedia Networks", *IEEE Wireless Communications*, pp. 76–83, October 2002.
6. R. Knopp and P.A. Humblet, "Information Capacity and Power Control in Single-Cell Multiuser Communications", in *Proceedings of IEEE International Conference on Communications (ICC'95)*, Seattle, USA, June 1995.
7. D. Wu and R. Negi, "Utilizing Multiuser Diversity for Efficient Support of Quality of Service Over a Fading Channel", in *Proceedings of IEEE ICC'03*, Anchorage, Alaska, USA, May 11–15, 2003.

Channel-Aware Earliest Deadline Due Fair Scheduling

8. A. Jalali, R. Padovani, and R. Pankaj, "Data Throughput of CDMA-HDR a High Efficiency-High Data Rate Personal Communication Wireless System", in *Proceedings of Vehicular Technology Conference 2000-spring*, vol. 3.
9. Y. Liu, S. Gruhl, and E. Knightly, "WCFQ: An Opportunistic Wireless Scheduler with Statistical Fairness Bounds", *IEEE Transactions on Wireless Communication*, September 2003.
10. Y. Liu and E. Knightly, "Opportunistic Fair Scheduling Over Multiple Wireless Channels", in *Proceedings of IEEE INFOCOM*, 2003.
11. X. Liu, E. K. P. Chong, and N. B. Shroff, "Transmission Scheduling for Efficient Wireless Utilization", in *Proceedings of IEEE INFOCOM*, 2000.
12. X. Liu, E. K. P. Chong, and N. B. Shroff, "Optimal Opportunistic Scheduling in Wireless Networks", in *Proceedings of IEEE INFOCOM*, 2000.
13. S. S. Kulkarni and C. Rosenberg, "Opportunistic Scheduling Policies for Wireless Systems with Short Term Fairness Constraints", in *Proceedings of IEEE GLOBECOM*, December 2003.
14. Z. Jiang and N. K. Shankaranarayana, "Channel Quality Dependent Scheduling for Flexible Wireless Resource Control", Globecom 2001, Texas.
15. P. Bender, P. Black, M. Grob, R. Padovani, N. Sindhushayana, and A. Viterbi, "CDMA/HDR: a Bandwidth Efficient High Speed Wireless Data Service for Nomadic Users", *IEEE Communications Magazine*, vol. 38, no. 7, pp. 70–77, July 2000.
16. 3rd Generation Partnership Project; Technical Specification Group Radio Access Network, "Physical Layer Aspects of UTRA High Speed Downlink Packet Access (Release 4) tr25.848 v4.0.0", Available: <http://www.3gpp.org>, Mar. 2001.
17. 3rd Generation Partnership Project; Technical Specification Group Radio Access Network, "Services Provided by the Physical Layer (Release 4) TS 25.302 v4.1.0", Available: <http://www.3gpp.org>, June 2001.
18. P. Bhagwat, A. Krishna, and S. Tripathi, "Enhancing Throughput Over Wireless LANs Using Channel State Dependent Packet Scheduling", in Proc. INFOCOM96, pp. 1133-1140, Mar. 1996.
19. S. Lu, V. Bharghavan, and R. Sirkant, "Fair Scheduling in Wireless Packet Networks", *IEEE Trans. on Networking*, pp. 473–89, Aug. 1999.
20. T. S. Ng, I. Stoica, and H. Zhang, "Packet Fair Queueing Algorithms for Wireless Networks with Location-Dependent Errors", in Proc. INFOCOM98, pp. 1103–1111, Mar. 1998.
21. R. Ramanathan and P. Agrawal, "Adapting Packet Fair Queueing Algorithms to Wireless Networks", in Proc. MOBICOM98, pp. 1–9, Oct. 1998.
22. S. Shakkottai and R. Srikant, "Scheduling Real-Time Traffic with Deadlines Over a Wireless Channel", in *Proceedings of ACM Workshop on Wireless and Mobile Multimedia*, Seattle, WA, August 1999.
23. A.L. Stolyar and K. Ramanan, "Largest Weighted Delay First Scheduling: Large Deviations and Optimality", *Annals of Applied Probability*, Vol. 11, No.1, pp. 1–48, 2001.
24. K. Chang and Y. Han, "QoS-Based Adaptive Scheduling for a Mixed Service in HDR System", *PIMRC 2002*.
25. S. Shakkottai and A. L. Stolyar, "Scheduling for Multiple Flows Sharing a Time-Varying Channel: The Exponential Rule", *Analytic Methods in Applied Probability*, Vol. 207, pp. 185–202, 2002.



Khaled M. F. Elsayed (S90–M95–SM02) received his B.Sc. (honors) in electrical engineering and M.Sc. in engineering mathematics from Cairo University in 1987 and 1990 respectively. He received his Ph.D. in computer science and computer engineering from North Carolina

K.M.F.Elsayed and A. K.F. Khattab

State University in 1995. He is now an Associate Professor in Cairo University, Egypt and is an independent telecommunications consultant. Between 1995 and 1997, he was a member of scientific staff with Nortel Wireless Systems Engineering in Richardson, TX.

Dr. Elsayed was the editor for the Internet technology series of the IEEE Communications Magazine from 1998 until 2002. He has served on technical program committees for several IEEE, IFIP, and ITC conferences. He was the technical co-chair for IFIP MWCN 2003 conference in Singapore. He also served as an expert evaluator for the European Commission FP5 and FP6 programmes. His research interest is in the area of performance evaluation of communication networks including IP, wireless and optical networks.



Ahmed Khattab received his B.Sc. (honors) and MS.C in Electronics and Communications Engineering from Cairo University in 2002 and 2004 respectively. Since August 2005, he is pursuing his PhD degree at Rice University, Texas. His research interests are in wireless networking and radio resource management.