

# Service Differentiation Without Prioritization in IEEE 802.11 WLANs

Suong H. Nguyen, Lachlan L. H. Andrew and Hai L. Vu  
Centre for Advanced Internet Architectures, Faculty of ICT  
Swinburne University, Hawthorn, VIC 3122, Australia  
{hsnguyen, landrew, hvu}@swin.edu.au

**Abstract**—Wireless LANs carry a mixture of traffic, with different delay and throughput requirements. The usual way to provide low-delay services is to give priority to such traffic. However this creates an incentive for throughput sensitive traffic also to use this service, which degrades overall network performance. We propose to allow applications to trade off delay for throughput, without giving preference to one class over another, by simultaneously scaling IEEE 802.11’s  $CW_{\min}$  and  $TXOP$  limit parameters.

We provide a model of this scheme with two traffic classes, and show that increasing  $CW_{\min}$  and  $TXOP$  limit in equal proportion reduces, but does not eliminate, the incentive for bulk data users to use the low-delay service. We show that subtracting a small constant from  $CW_{\min}$  eliminates this incentive, while still giving improved performance to both classes.

## I. INTRODUCTION

Wireless networks carry a diverse mix of traffic, from voice with tight delay constraints to bulk file downloads with only long-term throughput requirements. Efficient use of the network requires services tailored to each of these traffic classes. The traditional approach to providing quality of service (QoS) is to prioritize real-time traffic at the expense of data traffic, as done by the default parameter setting of the IEEE 802.11e Enhanced Distributed Channel Access (EDCA) standard [1]. This creates an incentive for data applications to use the class intended for real-time users to gain higher share of resources. This can degrade network performance drastically [2] and QoS differentiation no longer occurs when all data users use the highest priority class [3]. To cope with this, policing mechanisms have been proposed [4], which increase the complexity of the network.

As an alternative to prioritization, we propose a simple scheme that provides better service for both throughput- and delay-sensitive traffic, and encourages applications to use the service designed for them. The approach is based on the IEEE’s standard 802.11e [1] for carrier sense multiple access with collision avoidance (CSMA/CA) using exponential backoff. Note that unlike [5], [6], we are not combatting deliberately malicious users who violate the protocol, but rather addressing application writers who optimize their code based on measured performance using all the available services.

To support QoS, 802.11e uses EDCA, in which the access point (AP) selects four Access Categories (ACs) which stations can use. The ACs use different values of four MAC parameters:  $CW_{\min}$ ,  $CW_{\max}$ ,  $TXOP$  limit and  $AIFS$ . In particular,

TABLE I  
DEFAULT MAC PARAMETERS IN 802.11e EDCA

AC	$CW_{\min}$	$CW_{\max}$	AIFSN	$TXOP$ limit (for DSSS PHY)
Data (AC_BE)	31	1023	3	1 packet
Real-time (AC_VO)	7	15	2	3.264 ms

$CW_{\min}$  controls how long a station waits before transmission and  $TXOP$  limit controls how much it can transmit per channel access. Note that applications cannot choose arbitrary combinations of these parameters, but only those permitted by the AP.

The 802.11e EDCA standard recommends four particular combinations of parameters. The parameters intended for throughput-sensitive bulk data and delay-sensitive voice are shown in Table I, taken from Table 7-37 of [1]. In particular, the AC for real-time service is given uniformly higher priority than that for data service; every parameter is set to a more aggressive value. This priority-based QoS provision works well provided the high-priority class is only used by low-throughput real-time traffic. However, when applications are optimized based on performance measurements, the priority-based approach creates an incentive for all users to use the class AC\_VO, resulting in no QoS differentiation.

Instead, we seek to provide service differentiation without prioritizing one class over another, by choosing ACs such that some parameters are less aggressive whenever others are more aggressive. The aim of service differentiation without priority is to provide “different but fair” service for different types of traffic, by allowing users to choose different points on a throughput-delay tradeoff curve. Our proposal consists of choosing the combinations of parameters allowed by the AP. As such, it does not require any additional mechanisms such as fair queueing or traffic policing.

Fair differentiated service has already been proposed in [7]–[10] for wired networks. This has not been widely deployed, because it requires signals to be sent from the application to the core network. In contrast, for wireless links connected directly to the host running the application, no protocol changes are needed. Prior work in the wireless context can be divided into rewarding schemes [3], [11] and pricing schemes [2], [12]. The first approach uses 802.11e’s contention-free period

(CFP) to provide extra throughput to the data class, which is problematic because current wireless NICs do not implement the CFP. The pricing approach either requires micropayments of monetary prices, which makes implementation difficult, or must impose prices through some other form of service degradation such as packet drops, which seems counter-productive.

We will go through two steps to design our scheme.

First, in Section II, we propose a scheme called “proportional tradeoff” which provides better service for both traffic types, by judicious choice of  $CW_{\min}$  and  $TXOP$  limit for the two ACs. To gain insight into the scheme, we present a model of 802.11 WLANs in Section III with these two types of traffic from which the advantage of the scheme is shown in Section IV.

Second, we use a game theoretic model presented in Section III to investigate the incentives provided by the proportional tradeoff scheme. We show that our proposed scheme does not entirely eliminate data users’ incentive to use the real-time class, despite providing better service for data users than the default EDCA parameters. Hence the proposed scheme is modified in Section V by reducing the data class’s  $CW_{\min}$  slightly to give throughput-sensitive applications the incentive to use the bulk-data service. This modification is called “proportional incentive adjusted” or “PIA”. It is shown analytically that this gives the bulk-data class higher throughput than the real-time class under all traffic loads. Simulation shows that it also gives better delay performance to the real-time class than the parameter settings without service differentiation do.

## II. PROPOSED PROPORTIONAL TRADEOFF SCHEME

Our aim is to provide service differentiation for data traffic and real-time traffic, without the prioritization implied by standard 802.11e parameters. To this end, we now propose a mechanism which improves service for both types of traffic by increasing  $CW_{\min}$  and  $TXOP$  limit. We do not use the AIFS parameter because it provides load-dependent prioritization which makes it difficult to achieve a “fair” service differentiation.

In particular, we define two service classes: “bulk data” abbreviated by B and “real-time” abbreviated by R. Users of class B can transmit more packets per channel access but less often while users of class R transmit more often but only one packet per channel access. To achieve this, class B has higher  $TXOP$  limit but commensurately higher  $CW_{\min}$ . This is similar to the method used in [13] to ensure fairness. Let  $W_t$  and  $W_d$  be the values of  $CW_{\min}$  used by classes B and R respectively. Then

$$W_t = \eta W_d, \quad (1)$$

where  $\eta > 1$  is the number of packets sent by class B within its  $TXOP$  limit. Class B is designed for throughput-sensitive data users, and class R is for delay-sensitive real-time users.

The logic behind this scheme is that real-time traffic requires low delay and often has only one packet to send at a time but the packet needs to be sent as soon as possible; hence, it always uses class R. In contrast, a data source requires high

throughput; hence, it may be willing to wait a little longer, if this increases the amount it can transmit once it gains access to the channel.

We will show that when data users use class B, this scheme actually improves service for both types of traffic. This benefit comes from the reduction of collision probability in the network due to the lower attempt probability of data sources.

Note that MAC parameter setting in our scheme is per flow, not per node, because each EDCA station can support different flows with different types of traffic.

## III. MODEL

In this section, we present two models: wireless and game models. The wireless model studies the performance of our proportional tradeoff scheme given user choices. The game model studies the user incentive it induces. Note that these models apply to arbitrary MAC parameters, not only those satisfying (1). This allows them to be used to study the PIA scheme in Section V.

Consider a network consisting of  $N_u \geq 0$  non-saturated sources with arrival rate  $\lambda$ , and  $N_s \geq 1$  saturated sources, in an infrastructure topology. The non-saturated sources represent the real-time users while data users are modeled as saturated sources. For modeling simplicity, we assume each station transmits packets of only one source, although this is not required by the scheme itself. Each unsaturated station uses class R (sends a burst of one packet per channel access) while each saturated station can use either class R or class B (sends a burst of  $\eta \geq 1$  packets per channel access).

### A. Wireless model

We now present, without justification, the model derived in [14]. Note that this model has been validated extensively in [14], showing its accuracy under a wide range of scenarios.

The model uses the following assumptions:

- Channel conditions are ideal (no channel errors, hidden terminals or capture effect) and EDCA operates in basic mode (no RTS/CTS).
- Unsaturated stations always do backoff before a transmission. This standard assumption (see, e.g., [15], [16]) is reasonable in the presence of saturated stations.
- Sources perform binary exponential backoff until they successfully transmit a packet (there is no retransmission limit and  $CW_{\max} = \infty$ ). We do not claim that this assumption is realistic; however, simulations show that qualitative results from this model still hold when these two backoff parameters are truncated as in the standard.

The model uses the following notation:

- Subscripts  $s$ ,  $t$ ,  $d$ , and  $u$  denote any saturated user, a saturated user using class B, a saturated user using class R, and an unsaturated user using class R, respectively.
- $N_k$  ( $k \in \{s, t, d, u\}$ ) denotes the number of users of type  $k$ . Note that  $N_s = N_t + N_d$ .
- $W_k$  ( $k \in \{t, d, u\}$ ) is the minimum contention window of users of type  $k$ . Recall that a saturated user using R

transmits only one packet per channel access and has  $W_d = W_u$ .

- $\eta$  is the *TXOP limit* of a user of type  $t$ , which is the number of packets it can send per channel access<sup>1</sup>;  $l_{sat}$  and  $l_{nonsat}$  are the payload length of a packet from saturated and unsaturated sources, respectively. It is assumed that  $l_{nonsat} < l_{sat}$ .
- $\tau_k$  ( $k \in \{t, d, u\}$ ) is the probability that a user of type  $k$  attempts to transmit in a given slot.
- $p_k$  ( $k \in \{t, d, u\}$ ) is the collision probability of a packet from a user of type  $k$ .

The inputs to our model are  $N_t, N_d, N_u, W_t, W_u = W_d, \eta, \lambda, l_{sat}$  and  $l_{nonsat}$ . The outputs are attempt probabilities  $\tau_k$ , collision probabilities  $p_k$  ( $k \in \{t, d, u\}$ ), and the throughput of saturated sources.

Central to the model is a set of fixed-point equations, where the attempt probabilities of saturated and unsaturated sources are expressed in terms of the collision probabilities of saturated and unsaturated sources, and vice versa. The fixed point equations are as follows.

1) *Fixed point model*: First, the attempt probability of a saturated source of type  $k \in \{t, d\}$  is

$$\tau_k = \frac{\mathbb{E}[\text{Attempts per burst}]}{\mathbb{E}[\text{Slots per burst}]} = \frac{\sum_{j=0}^{\infty} p_k^j}{\sum_{j=0}^{\infty} (\mathbb{E}[U_{k,j}] + 1)p_k^j}$$

where  $U_{k,j}$  is a random variable representing the number of backoff slots in the  $j$ th backoff stage of a saturated station.  $U_{k,j}$  is an integer uniformly distributed on  $[0, 2^j W_k - 1]$ , whence

$$\mathbb{E}[U_{k,j}] = \frac{2^j W_k - 1}{2}.$$

Then,

$$\tau_k = \frac{2}{W_k \frac{1-p_k}{1-2p_k} + 1}, \quad k \in \{t, d\} \quad (2a)$$

Second, the attempt probability of an unsaturated source is

$$\begin{aligned} \tau_u &= \frac{\mathbb{E}[\text{Attempts per source per sec}]}{\mathbb{E}[\text{Number of system slots per sec}]} \\ &= \frac{\mathbb{E}[\text{Packets per source per sec}] \mathbb{E}[\text{Attempts per packet}]}{\mathbb{E}[\text{Number of system slots per sec}]} \\ &= \frac{\lambda \sum_{j=0}^{\infty} p_u^j}{(1/\mathbb{E}[Y])} \end{aligned}$$

where  $Y$  is the (random) virtual slot time, which is the time interval between the starts of two consecutive decrements of the backoff counter. Thus

$$\tau_u = \lambda \mathbb{E}[Y] \frac{1}{1-p_u}. \quad (2b)$$

Finally, the collision probability of a station is the probability that at least one of the remaining stations transmits in a given slot, namely

$$p_k = 1 - \frac{(1-\tau_t)^{N_s-N_d} (1-\tau_d)^{N_d} (1-\tau_u)^{N_u}}{1-\tau_k}, \quad k \in \{t, d, u\} \quad (2c)$$

<sup>1</sup>This differs slightly from [1] in which *TXOP limit* is a duration.

The fixed point is between the attempt probabilities in (2a)–(2b), and the collision probabilities in (2c), with the mean slot duration  $\mathbb{E}[Y]$  given by

$$\mathbb{E}[Y] = a_\sigma \sigma + a_u T_u + a_d T_d + a_t T_t \quad (3a)$$

$$a_\sigma = (1-\tau_t)^{N_t} (1-\tau_d)^{N_d} (1-\tau_u)^{N_u} \quad (3b)$$

$$a_u = (1-(1-\tau_u)^{N_u}) (1-\tau_t)^{N_t} (1-\tau_d)^{N_d} \quad (3c)$$

$$a_t = N_t \tau_t (1-\tau_t)^{N_t-1} (1-\tau_d)^{N_d} (1-\tau_u)^{N_u} \quad (3d)$$

$$a_d = 1 - (a_t + a_u + a_\sigma) \quad (3e)$$

where  $\sigma, T_u, T_d$ , and  $T_t$  are the duration of, respectively, idle slots, slots during which transmissions of only non-saturated stations occur, slots during which a collision involving at least one saturated station occurs or a successful transmission of a saturated station using class  $R$  occurs, and slots during which a successful transmission of a saturated station using class  $B$  occurs;  $a_\sigma, a_u, a_d$ , and  $a_t$  are the respective probabilities.

2) *Throughput of data users*: Two measures of throughput for saturated users are of interest here. In addition to the true throughput in packets/s, denoted  $S$ , some of our results apply to the more tractable measure of throughput in packets/slot, denoted  $C$ .

The throughput in packets/slot is given by

$$C_t = \tau_t (1-p_t) \eta \quad (4a)$$

$$C_d = \tau_d (1-p_d) \quad (4b)$$

where subscripts  $t$  and  $d$  denote saturated sources using class  $B$  and  $R$ , respectively.

The throughput in packets/s is given by

$$S_k = \frac{C_k}{\mathbb{E}[Y]}, \quad k \in \{t, d\}. \quad (5)$$

## B. Game model

1) *Description*: The above WLANs can be modeled as a game in which users are players. A player  $i$  can choose an action which is either to use class  $B$  or to use class  $R$ . Based on the actions of other players and its own action, player  $i$  will get a payoff, which is its throughput if  $i$  is a saturated user or the reciprocal of delay if  $i$  is an unsaturated station.

Using class  $R$  is a dominant strategy for unsaturated stations, since it reduces their delay regardless of what other stations do. For this reason, we will not treat unsaturated stations as players, but simply model their effect on the throughput obtained by the saturated users.

2) *Model*: A game of the wireless network described above is denoted by a quadruple  $\langle \mathcal{P}, (A_i)_{i \in \mathcal{P}}, (u_i)_{i \in \mathcal{P}}, N_u \rangle$  where

- $\mathcal{P} = \{1, \dots, N_s\}$ , the set of players, contains the saturated users.
- For every  $i \in \mathcal{P}$ ,  $A_i = \{B, R\}$  is the set of actions available to player  $i$ , where action  $B$  is to use MAC parameters  $(CW_{\min}, TXOP) = (W_t, \eta)$ , and action  $R$  is to use MAC parameters  $(W_u, 1)$ . Note that all the players have the same action space.
- For every  $i \in \mathcal{P}$ , the payoff  $u_i(a)$  is the throughput of player  $i$  under action  $a$  in which each player  $l$  plays  $a_l$ .

There are two forms of the game, corresponding to the two types of throughput described in Section III-A2:

- Game 1:  $u_i(a)$  is given by  $C_t$  of (4) if  $a_i = B$ , or  $C_d$  otherwise;
- Game 2:  $u_i(a)$  is given by  $S_t$  of (5) if  $a_i = B$ , or  $S_d$  otherwise.

Note that the values of  $C_t$ ,  $C_d$ ,  $S_t$  and  $S_d$  depend implicitly on the entire action profile  $a$ .

Note that this is a symmetric game [17], since each player has the same opportunities, and for each player, the same actions yield the same payoffs.

The properties of a game can be investigated using the wireless model (2)–(5).

#### IV. PROPERTIES OF PROPORTIONAL TRADEOFF SCHEME

In this section, we first use the MAC model (2)–(5) to derive some properties of the proportional tradeoff scheme. Then, we validate these results using ns-2 simulation.

Recall that under the proportional scheme, the relationship between the  $CW_{\min}$  of classes B and R is given by (1).

##### A. Theoretical results

We will investigate properties of the proportional scheme, both when data users always use class B (“simple” users), and where they use whichever class gives them higher throughput (“measurement driven” users). Even though our results are proved for a network with only data users and for unbounded  $CW_{\max}$  and number of retransmissions, we will show by simulation that they apply when these assumptions are lifted.

1) *Simple users:* Here we use the wireless model with  $N_d = 0$  to study network performance under the proposed scheme. The following theorem implies that our proposed scheme with  $\eta > 1$  provides better service for data users than does the case of no service differentiation ( $\eta = 1$ ).

*Theorem 1:* Under the wireless model based on (2)–(5) with  $W_t = \eta W_u$  ( $\eta \geq 1$ ,  $W_u > 4$ ),  $N_d = N_u = 0$ , the throughput in packets/s of each data user increases with  $\eta$ .

The above theorem, proved in [18], shows the advantage of the proposed scheme to data users. To show its benefit to real-time users, we will use simulation in Section IV-B.

2) *Measurement driven users:* When data users are measurement driven (using whichever class gives them higher throughput), we use the game model to analyze network performance. Under the proportional scheme, the action space of the game is

$$A_0 = \{(\eta W_u, \eta), (W_u, 1)\}, \quad \text{for a given } \eta > 1$$

where  $(\eta W_u, \eta)$  and  $(W_u, 1)$ , respectively, are MAC parameters of class B and R.

The following theorem states that, when the actions of all but one user are the same, the remaining user is better off by using class R. This implies that the proportional tradeoff scheme still creates an unfortunate incentive for a data user to choose class R, although this incentive is small. The theorem

TABLE II  
802.11B MAC AND PHY PARAMETERS

Parameter	Symbol	Value
Data bit rate	$r_{data}$	11 Mbps
Control bit rate	$r_{ctrl}$	1 Mbps
PHYS header	$T_{phys}$	192 $\mu s$
MAC header	$l_{mac}$	288 bits
UDP/IP header	$l_{udpip}$	160 bits
ACK packet	$l_{ack}$	112 bits
Slot time	$\sigma$	20 $\mu s$
SIFS	$T_{sifs}$	10 $\mu s$
Retry limit		7

TABLE III  
MAC PARAMETERS OF CLASSES B AND R

Class	$CW_{\min}$	$CW_{max}$	AIFSN	$TXOP$ limit (packets)
B	$W_t$	$2^5 W_t$	2	$\eta$
R	$W_u$	$2^5 W_u$	2	1

is proven in [18] with the following approximation of (2a) for  $W_k \gg 1$ :

$$\tau_k \approx \frac{2}{W_k} \frac{1 - 2p_k}{1 - p_k}. \quad (6)$$

*Theorem 2:* Consider the wireless model based on (2)–(4) with (2a) replaced by (6), in the game  $\langle \mathcal{P}, (A_0)_i \in \mathcal{P}, (C_i)_{i \in \mathcal{P}}, 0 \rangle$  with  $W_i > 4$ . For action profiles in which all users  $i \neq 1 \neq j$  have the same action  $a_i = a_j$ , data user 1 gets higher throughput per slot when using class R than when using class B.

Although the throughput in Theorem 2 is in packets/slot, simulation demonstrates that this result still holds for packets/s.

Recall that an action profile is a *Nash equilibrium* if no player can get higher payoff by changing its action while others keep theirs unchanged [19]. From Theorem 2, the action profile with all data users using class R is the unique symmetric Nash equilibrium. Then, according to Theorem 1, the throughput of each data user at this Nash equilibrium is less than that when all of data users use class B. Section V will consider an improved scheme that avoids that issue.

##### B. Simulation results and discussion

Recall that the properties of proportional tradeoff scheme in Sec. IV-A are proved for a network with only data users and for unbounded  $CW_{\max}$  and number of retransmissions. In this section, we will use simulation to validate those in more general scenarios with both data and real-time users, and a limited number of retransmissions. The simulations used ns-2.33 [20] with the EDCA package [21].

We simulated a network as described in Section III. All stations use the user datagram protocol (UDP). The general MAC and physical layer parameters are set to the default values in IEEE 802.11b, as shown in Table II.

The MAC parameters specific to classes B and R used in the proportional scheme are given in Table III with  $W_t = \eta W_u$ .

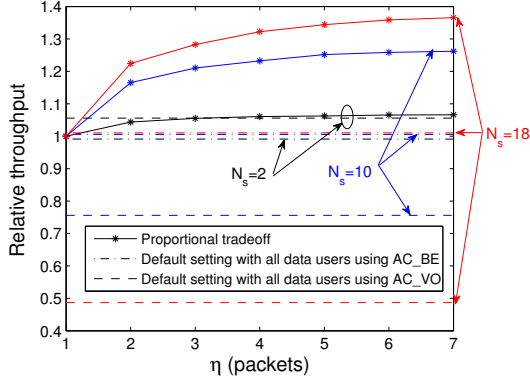


Fig. 1. Throughput of a data user as a function of class B's  $TxOP$  limit in packets ( $\eta$ ), scaled by that of the proportional scheme at  $\eta = 1$ . ( $\lambda = 20$  packets/s,  $l_{sat} = 1040$  bytes,  $l_{nonsat} = 100$  bytes,  $N_u = 4$ ,  $W_u = 32$ ,  $W_t = \eta W_u$ .)

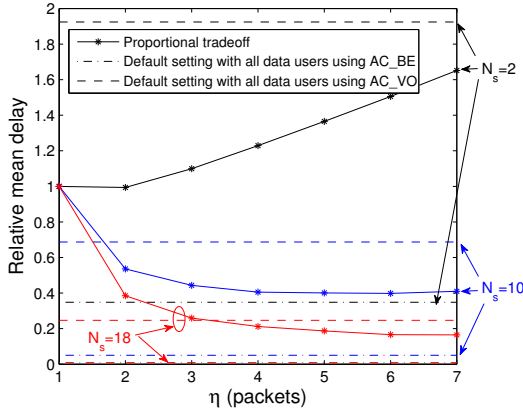


Fig. 2. Mean delay of a real-time user as a function of class B's  $TxOP$  limit in packets ( $\eta$ ), scaled by that of the proportional scheme at  $\eta = 1$ . ( $\lambda = 20$  packets/s,  $l_{sat} = 1040$  bytes,  $l_{nonsat} = 100$  bytes,  $N_u = 4$ ,  $W_u = 32$ ,  $W_t = \eta W_u$ .)

1) *Simple users* ( $N_d = 0$ ): To highlight the advantage of our scheme, we compare it with the default EDCA parameters in the same scenarios (same  $N_s$ ,  $N_u$ ,  $\lambda$ ,  $l_{sat}$ ,  $l_{nonsat}$ ).

The throughput of a data user, and the mean delay and loss probability of a real-time user, under the proportional scheme are shown in Figs. 1 and 2 and Table IV respectively, as functions of  $\eta$  for different  $N_s$ . Note that in Figs. 1 and 2, the performance metric at each  $\eta$  is normalized by that at  $\eta = 1$ . Moreover, the performance metrics under the default parameter setting (Table I) with all data users using class AC\_BE and real-time users using class AC\_VO and under the default parameter setting with all users using class AC\_VO are also shown for comparison. The performance metric of the default parameter setting in Figs. 1 and 2 is also normalized by that of the proportional scheme at  $\eta = 1$ . Note that the actual throughput and mean delay degrade as  $N_s$  increases; however, the relative performances improve as  $N_s$  increases.

Since the relative throughput is greater than 1 for  $\eta > 1$

TABLE IV  
LOSS PROBABILITY OF A REAL-TIME USER

	$N_s$			
	$\eta$	2	10	18
Proportional tradeoff	1	–	1.62e-4	8.17e-4
	2	–	1.9e-5	8.1e-5
	3	–	–	1.4e-5
	4	–	–	1.4e-5
	5	–	–	0.8e-5
	6	–	–	–
	7	–	–	–
Default setting with all data users using AC_BE		–	–	3.1e-5
Default setting with all data users using AC_VO		1.0e-4	7.37e-2	26.07e-2

(“–” denotes no loss found during simulation time)

as shown in Fig. 1, the proportional scheme with  $\eta > 1$  always provides better service for data users than the scheme with no service differentiation ( $\eta = 1$ ). This corroborates the result of Theorem 1. Also our scheme provides significantly better throughput for data users than does the default parameter setting, either with all data users using class AC\_BE or with all data users using class AC\_VO.

Note that the benefit of the proposed scheme increases with the contention level in the network. In Fig. 2, when the load is high enough, our scheme with  $\eta > 1$  provides significant improvement in terms of mean delay of real-time users in comparison with the case of no service differentiation ( $\eta = 1$ ). At light load (e.g.  $N_s = 2$ ), the improvement is negligible. This is acceptable because delay only becomes a problem at high load. Figure 2 also suggests that at each network load, there exists an optimal value of  $\eta$  at which mean delay is minimum (e.g.  $\eta = 2$  at  $N_s = 2$  and  $\eta = 6$  at  $N_s = 10$ ). This optimal  $\eta$  increases with the network load. Besides, Table IV shows that the loss probability of real-time traffic decreases with the increase of  $\eta$ .

Compared with the default parameter setting which prioritizes real-time traffic with all data users using class AC\_BE, we expect the performance will be worse for real-time users under the proportional scheme. This is seen in Fig. 2 and Table IV. However, compared with the default parameter setting with all data users using class AC\_VO, our proportional scheme provides much better service for real-time users. (The apparent exception for  $N_s = 18$  is due to the much higher loss rate under the default setting when all users use AC\_VO; see also Fig. 8.)

Although the optimal  $\eta$  in our proportional scheme depends on traffic load, the majority of the benefit for both data and real-time users is obtained at  $\eta = 2$ . Figs. 1 and 2 suggest that increasing  $\eta$  beyond 7 or 8 does not improve performance significantly, which is because the contention level does not decrease much further then.

2) *One measurement driven user* ( $N_d \leq 1$ ): To see if a data user has an incentive to use class R, we consider the case that  $N_s - 1$  data users use class B while the other chooses whether to use action  $R$  (outcome  $a_B$ ) or not (outcome  $a_A$ ). The ratio

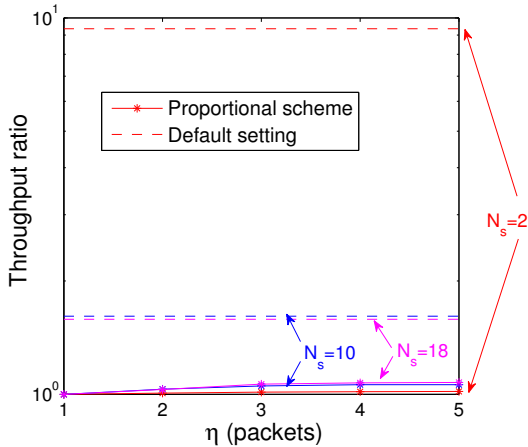


Fig. 3. Ratio of throughput in packets/s of a data user when it uses “real-time” class and “data” class while other data users use “data” class as a function of class B’s *TxOP limit* in packets ( $\eta$ ). ( $\lambda = 20$  packets/s,  $l_{sat} = 1040$  bytes,  $l_{nonsat} = 100$  bytes,  $N_u = 4$ ,  $W_u = 32$ ,  $W_t = \eta W_u$ .)

of throughput in packet/s of the undecided data user in  $a_B$  and  $a_A$  are shown in Fig. 3 as functions of  $\eta$  for different  $N_s$ . For comparison, we also show the ratio of throughput of a data user when it uses class AC\_VO and class AC\_BE while other data users use class AC\_BE under the default setting.

Observe in Fig. 3 that while other data users use “data” class, a data user can improve its throughput by using the real-time class under both the proportional scheme and the default parameter setting. However, the throughput improvement under the default parameter setting is much higher than under the proportional scheme. This shows that the incentive to use real-time class of data users under the proportional scheme is greatly reduced from the default parameter setting, even though it is not entirely eliminated.

3) *Many measurement driven users ( $N_d \leq N_s$ ):* The above shows that it is not a Nash equilibrium for all data users to use class B. To see if it is a Nash equilibrium for all to use class R, we consider the case that  $N_s - 1$  data users use class R while the other chooses whether to use action R (outcome  $a_N$ ) or not (outcome  $a_M$ ). The ratio of throughput  $S$  in packets/s of the undecided data user, user 1, in these two outcomes is shown in Fig. 4 as a function of  $\eta$  for different  $N_s$ .

Figure 4 shows that the throughput ratio is higher than 1 for  $\eta > 1$ , which means that while other data users use class R, a data user can improve its throughput by acting in the same way. This confirms that an action profile in which all data users choosing class R is a Nash equilibrium, as stated in Theorem 2.

However, this equilibrium gives a lower throughput than could be obtained when all data users use class B, as shown by the increase in relative throughput with  $\eta$  in Fig. 1. We next investigate a way to avoid this undesirable equilibrium.

## V. INCENTIVE ADJUSTED SCHEME

Section IV-B showed that for networks with both data and voice users, our proportional scheme can improve service for

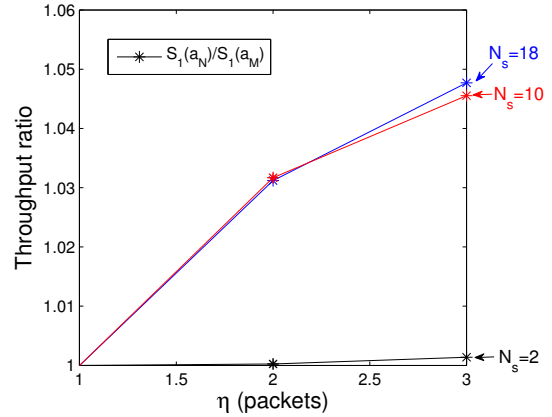


Fig. 4. Ratio of throughput in packets/s of the data user 1 when it uses class R ( $S_1(a_N)$ ) and class B ( $S_1(a_M)$ ) while other data users use class R as a function of class B’s *TxOP limit* in packets ( $\eta$ ). ( $\lambda = 20$  packets/s,  $l_{sat} = 1040$  bytes,  $l_{nonsat} = 100$  bytes,  $N_u = 4$ ,  $W_u = 32$ ,  $W_t = \eta W_u$ .)

both types of traffic relative to the scheme with no service differentiation, especially at high load. However, when a small fraction of data users chooses to use class R, their throughput can be slightly improved. Although the improvement is small, measurement-driven application design will still result in class R being chosen by throughput-sensitive applications. However, we will now show that our proposed scheme can easily be modified to eliminate this incentive issue. This is in contrast to priority-based schemes, which require explicit policing or pricing mechanisms.

### A. Description of the incentive adjusted scheme, PIA

We modify the proportional scheme by reducing  $CW_{\min}$  of class B by an amount  $\epsilon > 0$  and hence providing higher benefit for users of class B. The reduction in  $CW_{\min}$  for class B results in more throughput for a data user when it uses B compared with when it uses R, and thus data users have no incentive to use the real-time class. Recall that users can only select one of the access classes determined by the access point, and cannot choose arbitrary combinations of parameters.

Note that the performance of delay-sensitive users degrades as  $\epsilon$  increases, and so we would like to use the smallest  $\epsilon$  such that, for all network loads, bulk data users using class B get a higher throughput than those using class R. Any  $\epsilon$  larger than this will still induce the correct incentives. Finding the absolute smallest such  $\epsilon$  is beyond the scope of this paper, but it is shown in [18] that the model (2)–(5) implies that, for the typical case of  $W_u \geq 4$ , it is sufficient that  $\epsilon \geq \epsilon_0$ , where  $\epsilon_0$  is defined by

$$\epsilon_0 \equiv 4(\eta - 1).$$

Specifically, under the incentive adjusted scheme PIA, the action space in the game model has the form

$$A1 = \{(\eta W_u - \epsilon_0, \eta), (W_u, 1)\}, \quad \text{for a given } \eta > 1$$

where  $(\eta W_u - \epsilon_0, \eta)$  and  $(W_u, 1)$ , respectively, are MAC parameters of class B and R.

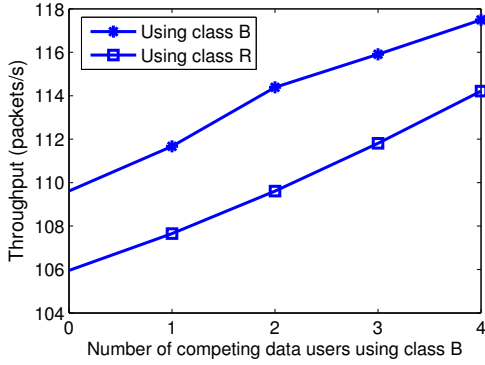


Fig. 5. Throughput of class-B and class-R data users as a function of the number of class-R data users  $N_d$ . ( $\lambda = 30$  pkts/s,  $l_{sat} = 1040$ B,  $l_{nonsat} = 100$ B,  $N_u = 6$ ,  $N_s = 5$ ,  $W_u = 32$ ,  $W_t = \eta W_u - \epsilon_0$ ,  $\eta = 2$ .)

### B. Evaluation of the incentive adjusted scheme

To highlight the advantage of the incentive adjusted scheme PIA when all data users are measurement driven, we compare it with the case of no service differentiation ( $\eta = 1$ ) and the default EDCA parameter setting (Table I). Recall that under PIA, data users have no incentive to use class R. In contrast, all data users have incentive to use the highest priority class in the default priority scheme.

In this section, the MAC parameters specific to classes B and R in our PIA scheme are also given in Table III with  $W_t = \eta W_u - \epsilon_0$ .

1) *Incentive compatibility*: We will first verify that throughput-sensitive users have an incentive to choose class B, regardless of the class chosen by other users. To do this, we simulated a scenario with  $N_s = 5$ , and plotted results in Fig. 5.

Figure 5 shows the throughput a saturated user obtains by using either class B or class R, as a function of the number of competing data users who use class B. Regardless of the choices of the other users, a given saturated user gets a higher throughput by choosing the bulk-data class. If all data users choose class B, then the total throughput is maximal.

This was the objective in selecting  $\epsilon_0$ . However, the result is stronger than is ensured by the theory in [18], since the latter only applied to the cases when either none or all of the competing data users used class R. Thus the simulation suggests that PIA is actually incentive compatible, and causes bulk data users to choose class B.

2) *Comparison with the default EDCA parameters*: We can now compare the performance of the proposed PIA with that of the default QoS classes, under the assumption that all users will use whatever class gives them the best performance. Under the default parameters, all users will use AC\_VO, and under PIA bulk data users will use class B and real-time users will use class R.

We look at the scenarios where  $\lambda = 20$  packets/s,  $l_{sat} = 1040$ B,  $l_{nonsat} = 100$ B,  $N_u = 4$ ,  $N_s = \{2, 10, 18\}$ ,  $W_u = 32$ ,  $W_t = \eta W_u - \epsilon_0$ ,  $\epsilon_0 = 4(\eta - 1)$  and  $\eta$  is varied. The case of no service differentiation ( $\eta = 1$ ) is included for comparison.

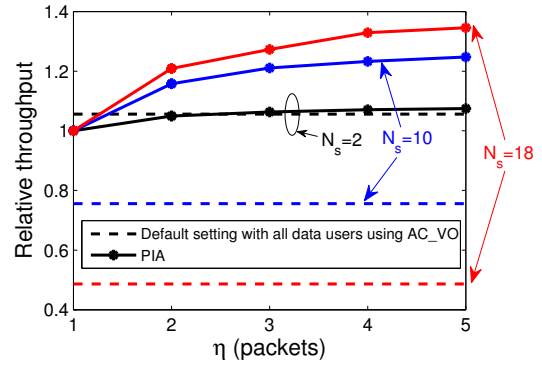


Fig. 6. Throughput of a data user as a function of class B's  $TxOP$  limit in packets ( $\eta$ ), scaled by that of PIA at  $\eta = 1$ . ( $\lambda = 20$  packets/s,  $l_{sat} = 1040$  bytes,  $l_{nonsat} = 100$  bytes,  $N_u = 4$ ,  $W_u = 32$ ,  $W_t = \eta W_u - \epsilon_0$  for  $\eta = \{1, 2, 3, 4, 5\}$ .)

The relative throughput of a saturated user under PIA is shown in Fig. 6 as functions of  $\eta$  for different numbers of saturated users,  $N_s$ . For comparison, the throughput under the default parameter setting (Table I) is also shown. The throughput is again normalized by that of PIA at  $\eta = 1$ .

The throughput increases faster with  $\eta$  under PIA than it did under the proportional scheme, which reflects the reduction in  $CW_{min}$ . This shows that PIA provides better service for data users than the proportional scheme without service differentiation ( $\eta = 1$ ), especially at high load. This is in contrast to the default parameter setting with all data users using real-time class (AC\_VO), for which the performance degrades rapidly at high load. For low load ( $N_s = 2$ ), the default parameter setting performs better than the system without QoS differentiation ( $\eta = 1$ ), because the more aggressive choice of  $CW_{min}$  is better matched to a small number of stations.

This improvement in throughput of PIA comes at the expense of increased delay for real-time users. Figs. 7 and 8 show the probability that a packet of a real-time user is successfully transmitted before a given delay, for different  $\eta$  and loads  $N_s = 2$  and  $N_s = 18$ .

Figure 7 shows that PIA at both  $\eta = 2$  and  $\eta = 5$  gives a higher probability that a packet is successfully delivered at a given delay than the default parameter setting with all data users using the real-time class. This means that the average packet delay under PIA is smaller. In this lightly loaded case,  $\eta = 2$  provides comparable service to  $\eta = 1$  (no service differentiation), and  $\eta = 5$  provides slight degradation, but less than that caused by the default prioritization.

In the heavily loaded case of Fig. 8, the cumulative distribution of delay for the default parameter setting never reaches 1, which indicates a high loss rate. In contrast, PIA has a low loss rate for all values of  $\eta$  tested, although some packets have very high delays. In this case, the benefit increases as  $\eta$  increases. Together with the result in Fig. 7, this implies that under PIA, the optimal  $\eta$  for real-time users increases with traffic load, as was observed for the proportional scheme. However, even using  $\eta = 2$  for all loads appears to provide improvement over

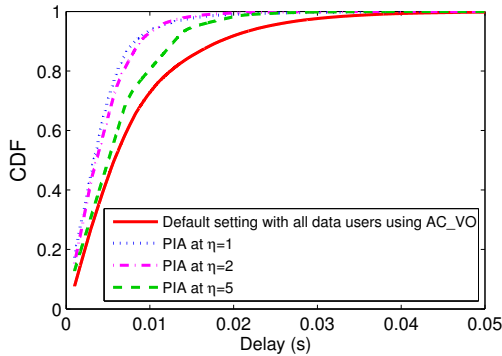


Fig. 7. CDF of delay: the probability a packet of real-time users is successfully delivered within a given time. ( $\lambda = 20$  pkts/s,  $l_{sat} = 1040B$ ,  $l_{nonsat} = 100B$ ,  $N_u = 4$ ,  $N_s = 2$ ,  $W_u = 32$ ,  $W_t = \eta W_u - \epsilon_0$  for  $\eta = \{1, 2, 5\}$ .)

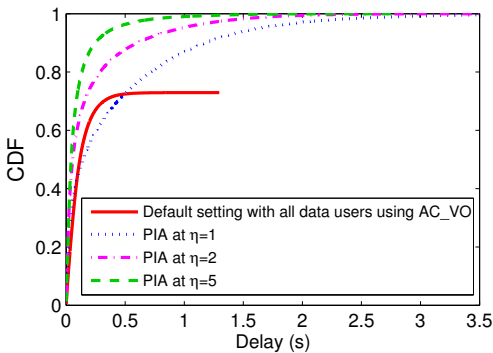


Fig. 8. CDF of delay: the probability a packet of real-time users is successfully delivered within a given time. ( $\lambda = 20$  pkts/s,  $l_{sat} = 1040B$ ,  $l_{nonsat} = 100B$ ,  $N_u = 4$ ,  $N_s = 18$ ,  $W_u = 32$ ,  $W_t = \eta W_u - \epsilon_0$  for  $\eta = \{1, 2, 5\}$ .) The delay under the default setting never exceeds 1.4s due to the finite retransmission limit and small  $CW_{min}$ .

the default parameters.

In summary, although the optimal  $\eta$  in PIA depends on traffic load, it is clear that when  $\eta = 2$ , PIA provides better service for both types of traffic under typical scenarios. This implies that when designing a network with an unknown number of users, PIA can be implemented by simply setting  $\eta = 2$  and  $\epsilon_0 = 4$ . Adaptive schemes that set  $\eta$  dependent on the estimated load are possible, but out of scope of this paper.

## VI. CONCLUSION

It is important to be able to provide differentiated services, without giving all users the incentive to use a “highest priority” class. This paper has shown through both analysis and simulation that allowing users to adjust  $CW_{min}$  and  $TXOP$  limit in the same proportion provides service differentiation in WLANs. This scheme improves service for both data and real-time traffic, especially at high load. However, it still provides a slight incentive for data users to use real-time class’s parameters. This misalignment of incentives can be removed by increasing  $CW_{min}$  by a slightly smaller factor than the  $TXOP$  limit. Our incentive adjusted scheme has many advantages over

prior proposals: it improves service for both data and real-time traffic and provides the correct incentives for application optimizers, while allowing easy implementation: a single set of 802.11e MAC parameters provides good performance over wide range of loads.

## ACKNOWLEDGMENT

This work was supported by Australian Research Council grants DP1095103 and FT0991594.

## REFERENCES

- [1] *Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications*, IEEE Standard 802.11<sup>TM</sup>, 2007.
- [2] M. H. Cheung, A. H. Mohsenian-Rad, V. W. S. Wong, and R. Schober, “Random Access Protocols for WLANs based on mechanism design,” in *Proc. IEEE Int. Conf. Commun.*, pp. 4892–4897, 2009.
- [3] J. Price, P. Nuggehalli, and T. Javidi, “Incentive Compatible MAC-Layer QoS Design,” in *Proc. CCNC 2008*.
- [4] C. Chun-Ting, S. N. Shankar and K. G. Shin, “Achieving per-stream QoS with distributed airtime allocation and admission control in IEEE 802.11e wireless LANs” in *Proc. INFOCOM* pp. 1584–1595, 2005.
- [5] M. Cagalj, S. Ganeriwal, I. Aad, and J.P. Hubaux, “On selfish behavior in CSMA/CA networks,” in *Proc. IEEE INFOCOM*, 2005.
- [6] L. Galluccio, “A Game-Theoretic Approach to Prioritized Transmission in Wireless CSMA/CA Networks,” in IEEE Vehicular Technology Conference, April 2009.
- [7] P. Hurley and J.Y. Le Boudec, “ABE: Providing a low-delay service within best effort,” *IEEE Network*, vol. 15, no. 3, pp. 60–69, 2001.
- [8] M. Karsten, Y. Lin, and K. Larson, “Incentive-Compatible Differentiated Scheduling,” in *Proc. HotNets IV*, 2005.
- [9] B. Gaidioz and P. Primet, “EDS: A new scalable service differentiation architecture for internet,” in *Proc. IEEE ISCC*, pp. 777–782, July 2002.
- [10] V. Firoiu, X. Zhang, and Y. Guo, “Best Effort Differentiated Services: Trade-off Service Differentiation for Elastic Applications,” in *Proc. IEEE ICT*, June 2001.
- [11] P. Nuggehalli, M. Sarkar, K. Kulkarni, and R. R. Rao, “Game-theoretic Analysis of QoS in Wireless MAC,” in *Proc. INFOCOM 2008*.
- [12] J. Price, P. Nuggehalli, and T. Javidi, “Pricing and QoS in Wireless Random Access Networks,” in *Proc. GLOBECOM 2008*.
- [13] N. H. Vaidya, A. Dugar, S. Gupta and P. Bahl, “Distributed fair scheduling in a wireless LAN,” *IEEE Trans. Mobile Computing*, vol. 4, no. 6, pp. 616–629, Nov/Dec 2005.
- [14] S. H. Nguyen, H. L. Vu and L. L. H. Andrew, “Performance analysis of 802.11e with saturated and non-saturated sources,” (Online). Available <http://caia.swin.edu.au/reports/091013A/CAIA-TR-091013A.pdf>.
- [15] X. Ling, L. X. Cai, J. W. Mark, and X. Shen, “Performance Analysis of IEEE 802.11 DCF with heterogeneous traffic,” in *Proc. IEEE Consumer Commun. Netw. Conf. (CCNC 2007)*, art. no. 4199105, pp. 49–53, 2007.
- [16] B. Bellalta, C. Cano, M. Oliver, and M. Meo, “Modeling the IEEE 802.11e EDCA for MAC parameter optimization,” in *Proc. IEEE Consumer Commun. Netw. Conf. (CCNC)*, vol. 1, pp. 390–394, 2006.
- [17] P. K. Dutta, *Strategies and games: theory and practice*. The MIT Press, 1999.
- [18] S. H. Nguyen, H. L. Vu and L. L. H. Andrew, “Service differentiation without prioritization in 802.11 WLANs,” Available at <http://caia.swin.edu.au/cv/snguyen/J2.pdf>.
- [19] D. Fudenberg and J. Tirole, *Game Theory*. The MIT Press, 1991.
- [20] “The network simulator ns-2,” Available at <http://www.isi.edu/nsnam/ns/>.
- [21] S. Wietholter and C. Hoene, “An IEEE 802.11e EDCF and CFB simulation model for ns-2,” Available at [http://www.tkn.tu-berlin.de/research/802.11e\\_ns2/](http://www.tkn.tu-berlin.de/research/802.11e_ns2/).