

Buffer Sizing for 802.11 Based Networks

Tianji Li, Douglas Leith, David Malone

Hamilton Institute, National University of Ireland Maynooth, Ireland

Email: {tianji.li, doug.leith, david.malone}@nuim.ie

Abstract—We consider the sizing of network buffers in 802.11 based networks. Wireless networks face a number of fundamental issues that do not arise in wired networks. We demonstrate that the use of fixed size buffers in 802.11 networks inevitably leads to either undesirable channel under-utilization or unnecessary high delays. We present two novel dynamic buffer sizing algorithms that achieve high throughput while maintaining low delay across a wide range of network conditions. Experimental measurements demonstrate the utility of the proposed algorithms in a production WLAN and a lab testbed.

Index Terms—IEEE 802.11, IEEE 802.11e, Wireless LANs (WLANs), Medium access control (MAC), Transmission control protocol (TCP), Buffer Sizing, Stability Analysis.

I. INTRODUCTION

In communication networks, buffers are used to accommodate short-term packet bursts so as to mitigate packet drops and to maintain high link efficiency. Packets are queued if too many packets arrive in a sufficiently short interval of time during which a network device lacks the capacity to process all of them immediately.

For wired routers, the sizing of buffers is an active research topic ([31] [5] [27] [32] [9]). The classical rule of thumb for sizing wired buffers is to set buffer sizes to be the product of the *bandwidth* and the average *delay* of the flows utilizing this link, namely the *Bandwidth-Delay Product* (BDP) rule [31]. See Section VII for discussion of other related work.

Surprisingly, however the sizing of buffers in wireless networks (especially those based on 802.11/802.11e) appears to have received very little attention within the networking community. Exceptions include the recent work in [21] relating to buffer sizing for voice traffic in 802.11e [2] WLANs, work in [23] which considers the impact of buffer sizing on TCP upload/download fairness, and work in [29] which is related to 802.11e parameter settings.

Buffers play a key role in 802.11/802.11e wireless networks. To illustrate this, we present measurements from the production WLAN of the Hamilton Institute, which show that the current state of the art which makes use of fixed size buffers, can easily lead to poor performance. The topology of this WLAN is shown in Fig. 23. See the Appendix for further details of the configuration used. We recorded RTTs before and after one wireless station started to download a 37MByte file from a web-site. Before starting the download, we pinged the access point (AP) from a laptop 5 times, each time sending 100 ping packets. The RTTs reported by the ping program was between 2.6-3.2 ms. However, after starting

the download and allowing it to continue for a while (to let the congestion control algorithm of TCP probe for the available bandwidth), the RTTs to the AP hugely increased to **2900-3400 ms**. During the test, normal services such as web browsing experienced obvious pauses/lags on wireless stations using the network. Closer inspection revealed that the buffer occupancy at the AP exceeded 200 packets most of the time and reached 250 packets from time to time during the test. Note that the increase in measured RTT could be almost entirely attributed to the resulting queuing delay at the AP, and indicates that a more sophisticated approach to buffer sizing is required. Indeed, using the A* algorithm proposed in this paper, the RTTs observed when repeating the same experiment fall to only **90-130 ms**. This reduction in delay does not come at the cost of reduced throughput, i.e., the measured throughput with the A* algorithm and the default buffers is similar.

In this paper, we consider the sizing of buffers in 802.11/802.11e ([1] [2]) based WLANs. We focus on single-hop WLANs since these are rapidly becoming ubiquitous as the last hop on home and office networks as well as in so-called “hot spots” in airports and hotels, but note that the proposed schemes can be easily applied in multi-hop wireless networks. Our main focus in this paper is on TCP traffic since this continues to constitute the bulk of traffic in modern networks (80–90% [35] of current Internet traffic and also of WLAN traffic [28]), although we extend consideration to UDP traffic at various points during the discussion and also during our experimental tests.

Compared to sizing buffers in wired routers, a number of fundamental new issues arise when considering 802.11-based networks. Firstly, unlike wired networks, wireless transmissions are inherently broadcast in nature which leads to the packet service times at different stations in a WLAN being strongly coupled. For example, the basic 802.11 DCF ensures that the wireless stations in a WLAN win a roughly equal number of transmission opportunities [19], hence, the mean packet service time at a station is an order of magnitude longer when 10 other stations are active than when only a single station is active. Consequently, the buffering requirements at each station would also differ, depending on the number of other active stations in the WLAN. In addition to variations in the mean service time, the distribution of packet service times is also strongly dependent on the WLAN offered load. This directly affects the burstiness of transmissions and so buffering requirements (see Section III for details). Secondly, wireless stations dynamically adjust the physical transmission rate/modulation used in order to regulate non-congestive channel losses. This rate adaptation, whereby the transmit rate may change by a factor of 50 or more (e.g. from 1Mbps to

54Mbps in 802.11a/g), may induce large and rapid variations in required buffer sizes. Thirdly, the ongoing 802.11n standards process proposes to improve throughput efficiency by the use of large frames formed by aggregation of multiple packets ([3] [18]). This acts to couple throughput efficiency and buffer sizing in a new way since the latter directly affects the availability of sufficient packets for aggregation into large frames.

It follows from these observations that, amongst other things, there does not exist a fixed buffer size which can be used for sizing buffers in WLANs. This leads naturally to consideration of dynamic buffer sizing strategies that adapt to changing conditions. In this paper we demonstrate the major performance costs associated with the use of fixed buffer sizes in 802.11 WLANs (Section III) and present two novel dynamic buffer sizing algorithms (Sections IV and V) that achieve significant performance gains. The stability of the feedback loop induced by the adaptation is analyzed, including when cascaded with the feedback loop created by TCP congestion control action. The proposed dynamic buffer sizing algorithms are computationally cheap and suited to implementation on standard hardware. Indeed, we have implemented the algorithms in both the NS-2 simulator and the Linux MadWifi driver [4]. In this paper, in addition to extensive simulation results we also present experimental measurements demonstrating the utility of the proposed algorithms in a testbed located in office environment and with realistic traffic. This latter includes a mix of TCP and UDP traffic, a mix of uploads and downloads, and a mix of connection sizes.

The remainder of the paper is organized as follows. Section II introduces the background of this work. In Section III simulation results with fixed size buffers are reported to further motivate this work. The proposed algorithms are then detailed in Sections IV and V. Experiment details are presented in Section VI. After introducing related work in Section VII, we summarize our conclusions in Section VIII.

II. PRELIMINARIES

A. IEEE 802.11 DCF

IEEE 802.11a/b/g WLANs all share a common MAC algorithm called the Distributed Coordinated Function (DCF) which is a CSMA/CA based algorithm. On detecting the wireless medium to be idle for a period $DIFS$, each wireless station initializes a backoff counter to a random number selected uniformly from the interval $[0, CW-1]$ where CW is the contention window. Time is slotted and the backoff counter is decremented each slot that the medium is idle. An important feature is that the countdown halts when the medium is detected busy and only resumes after the medium is idle again for a period $DIFS$. On the counter reaching zero, a station transmits a packet. If a collision occurs (two or more stations transmit simultaneously), CW is doubled and the process repeated. On a successful transmission, CW is reset to the value CW_{min} and a new countdown starts.

B. IEEE 802.11e EDCA

The 802.11e standard extends the DCF algorithm (yielding the EDCA) by allowing the adjustment of MAC parameters

T_{SIFS} (μs)	10
Idle slot duration (σ) (μs)	9
Retry limit	11
Packet size (bytes)	1000
PHY data rate (Mbps)	54
PHY basic rate (Mbps)	6
PLCP rate (Mbps)	6

TABLE I
MAC/PHY PARAMETERS USED IN SIMULATIONS, CORRESPONDING TO 802.11G.

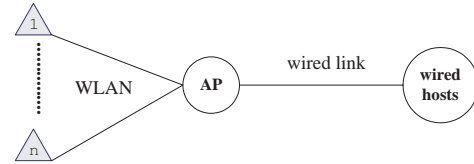


Fig. 1. WLAN topology used in simulations. Wired backhaul link bandwidth 100Mbps. MAC parameters of the WLAN are listed in Table I.

that were previously fixed. In particular, the values of $DIFS$ (called $AIFS$ in 802.11e) and CW_{min} may be set on a per class basis for each station. While the full 802.11e standard is not implemented in current commodity hardware, the EDCA extensions have been widely implemented for some years.

C. Unfairness among TCP Flow

Consider a WLAN consisting of n client stations each carrying one TCP upload flow. The TCP ACKs are transmitted by the wireless AP. In this case TCP ACK packets can be easily queued/dropped due to the fact that the basic 802.11 DCF ensures that stations win a roughly equal number of transmission opportunities. Namely, while the data packets for the n flows have an aggregate $n/(n+1)$ share of the transmission opportunities the TCP ACKs for the n flows have only a $1/(n+1)$ share. Issues of this sort are known to lead to significant unfairness amongst TCP flows but can be readily resolved using 802.11e functionality by treating TCP ACKs as a separate traffic class which is assigned higher priority [15]. With regard to throughput efficiency, the algorithms in this paper perform similarly when the DCF is used and when TCP ACKs are prioritized using the EDCA as in [15]. Per flow behavior does, of course, differ due to the inherent unfairness in the DCF and we therefore mainly present results using the EDCA to avoid flow-level unfairness.

D. Simulation Topology

In Sections III, IV and V-G, we use the simulation topology shown in Fig. 1 where the AP acts as a wireless router between the WLAN and the Internet. Upload flows originate from stations in the WLAN on the left and are destined to wired host(s) in the wired network on the right. Download flows are from the wired host(s) to stations in the WLAN. We ignore differences in wired bandwidth and delay from the AP to the wired hosts which can cause TCP unfairness issues on the wired side (an orthogonal issue) by using the same wired-part RTT for all flows. Unless otherwise stated, we use the

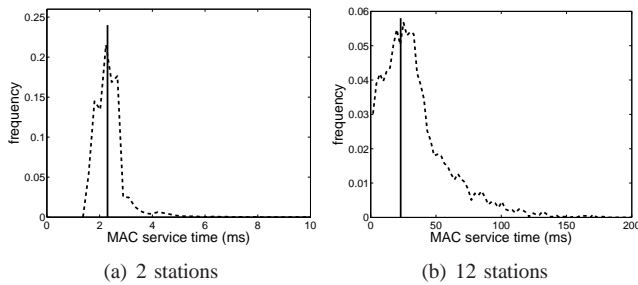


Fig. 2. Measured distribution of per packet MAC service time. Solid vertical lines mark the mean values of distributions. Physical layer data/basic rates are 11/1 Mbps.

IEEE 802.11g PHY parameters shown in Table I and the wired backhaul link bandwidth is 100Mbps with RTT 200ms. For TCP traffic, the widely deployed TCP Reno with SACK extension is used. The advertised window size is set to be 4096 packets (each has a payload of 1000 bytes) which is the default size of current Linux kernels. The maximum value of the TCP smoothed RTT measurements (sRTT) is used as the measure of the delay experienced by a flow.

III. MOTIVATION AND OBJECTIVES

Wireless communication in 802.11 networks is *time-varying* in nature, i.e., the mean service time and the distribution of service time at a wireless station vary in time. The variations are primarily due to (i) changes in the number of active wireless stations and their load (i.e. offered load on the WLAN) and (ii) changes in the physical transmit rate used (i.e. in response to changing radio channel conditions). In the latter case, it is straightforward to see that the service time can be easily increased/decreased using low/high physical layer rates. To see the impact of offered load on the service time at a station, Fig. 2 plots the measured distribution of the MAC layer service time when there are 2 and 12 stations active. It can be seen that the mean service time changes by over an order of magnitude as the number of stations varies. Observe also from these measured distributions that there are significant fluctuations in the service time for a given fixed load. This is a direct consequence of the stochastic nature of the CSMA/CA contention mechanism used by the 802.11/802.11e MAC.

This time-varying nature directly affects buffering requirements. Figure 3 plots link utilization¹ and max sRTT (propagation plus smoothed queuing delay) vs buffer size for a range of WLAN offered loads and physical transmit rates. We can make a number of observations.

First, it can be seen that as the physical layer transmit rate is varied from 1Mbps to 216Mbps, the minimum buffer size to ensure at least 90% throughput efficiency varies from about 20 packets to about 800 packets. No compromise buffer size exists that ensures both high efficiency and low delay across this range of transmit rates. For example, a buffer size of 80 packets leads to RTTs exceeding 500ms (even when only a single station is active and so there are no competing

wireless stations) at 1Mbps and throughput efficiency below 50% at 216Mbps. Note that the transmit rates in currently available draft 802.11n equipment already exceed 216Mbps (e.g. 300Mbps is supported by current Atheros chipsets) and the trend is towards still higher transmit rates. Even across the restricted range of transmit rates 1Mbps to 54Mbps supported by 802.11a/b/g, it can be seen that a buffer size of 50 packets is required to ensure throughput efficiency above 80% yet this buffer size induces delays exceeding 1000 and 3000 ms at transmit rates of 11 and 1 Mbps, respectively.

Second, delay is strongly dependent on the traffic load and the physical rates. For example, as the number of competing stations (marked as “uploads” in the figure) is varied from 0 to 10, for a buffer size of 20 packets and physical transmit rate of 1Mbps the delay varies from 300ms to over 2000ms. This reflects that the 802.11 MAC allocates available transmission opportunities equally on average amongst the wireless stations, and so the mean service time (and thus delay) increases with the number of stations. In contrast, at 216Mbps the delay remains below 500ms for buffer sizes up to 1600 packets.

Our key conclusion from these observations is that there exists no fixed buffer size capable of ensuring both high throughput efficiency and reasonable delay across the range of physical rates and offered loads experienced by modern WLANs. Any fixed choice of buffer size necessarily carries the cost of significantly reduced throughput efficiency and/or excessive queuing delays.

This leads naturally therefore to the consideration of adaptive approaches to buffer sizing, which dynamically adjust the buffer size in response to changing network conditions to ensure both high utilization of the wireless link while avoiding unnecessarily long queuing delays.

IV. EMULATING BDP

We begin by considering a simple adaptive algorithm based on the classical BDP rule. Although this algorithm cannot take advantage of statistical multiplexing opportunities, it is of interest both for its simplicity and because it will play a role in the more sophisticated A^* algorithm developed in the next section.

As noted previously, and in contrast to wired networks, in 802.11 WLANs the mean service time is generally time-varying (dependent on WLAN load and the physical transmit rate selected by a station). Consequently, there does not exist a fixed BDP value. However, we note that a wireless station can measure its own packet service times by direct observation, i.e., by recording the time between a packet arriving at the head of the network interface queue t_s and being successfully transmitted t_e (which is indicated by receiving correctly the corresponding MAC ACK). Note that this measurement can be readily implemented in real devices, e.g. by asking the hardware to raise an interrupt on receipt of a MAC ACK, and incurs only a minor computational burden. Averaging these per packet service times yields the mean service time T_{serv} . To accommodate the time-varying nature of the mean service time, this average can be taken over a sliding window. In this paper, we consider the use of exponential smoothing

¹Here the AP throughput percentage is the ratio between the actual throughput achieved using buffer sizes shown on the x-axis and the maximum throughput using the buffer sizes shown on the x-axis.

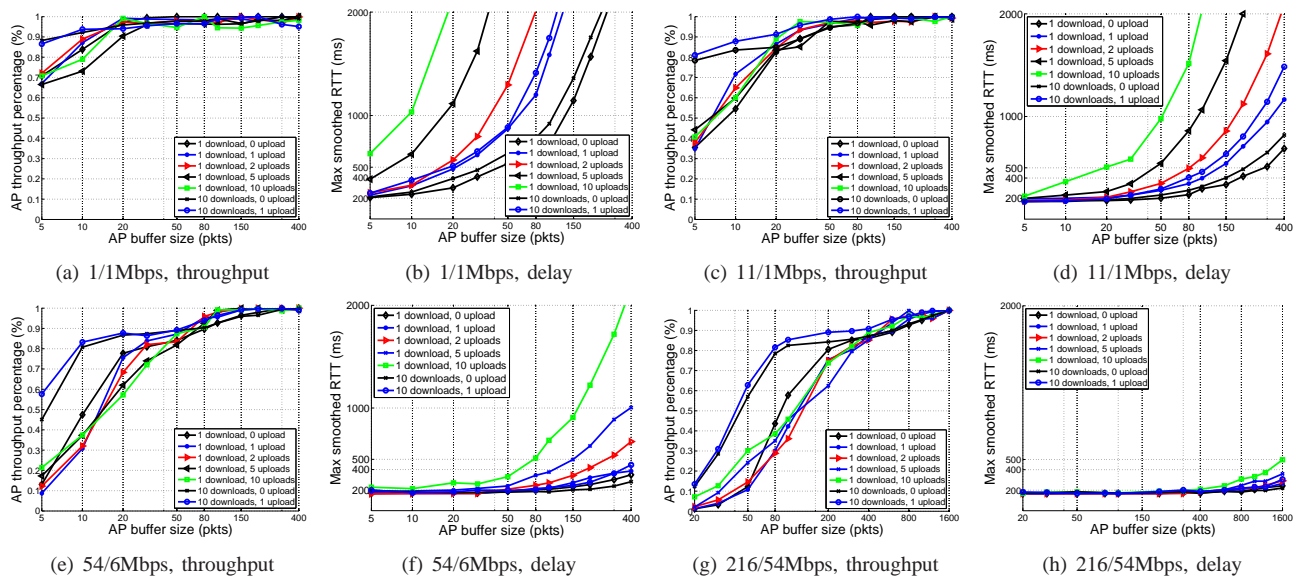


Fig. 3. Throughput efficiency and maximum smoothed round trip delays (max sRTT) for the topology in Fig. 1 when fixed size buffers are used. Here, the AP throughput efficiency is the ratio between the download throughput achieved using buffer sizes indicated on the x-axis and the maximum download throughput achieved using fixed size buffers. Rates before and after the '/' are used physical layer data and basic rates. For the 216Mbps data, 8 packets are aggregated into each frame at the MAC layer to improve throughput efficiency in an 802.11n-like scheme. The wired RTT is 200 ms.

$T_{serv}(k+1) = (1-W)T_{serv}(k) + W(t_e - t_s)$ to calculate a running average since this has the merit of simplicity and statistical robustness (by central limit arguments). The choice of smoothing parameter W involves a trade-off between accommodating time variations and ensuring the accuracy of the estimate – this choice is considered in detail later.

Given an online measurement of the mean service time T_{serv} , the classical BDP rule yields the following eBDP buffer sizing strategy. Let T_{max} be the target maximum queuing delay. Noting that $1/T_{serv}$ is the mean service rate, we select buffer size Q_{eBDP} according to $Q_{eBDP} = \min(T_{max}/T_{serv}, Q_{max}^{eBDP})$ where Q_{max}^{eBDP} is the upper limit on buffer size. This effectively regulates the buffer size to equal the current mean BDP. The buffer size decreases when the service rate falls and increases when the service rate rises, so as to maintain an approximately constant queuing delay of T_{max} seconds. We may measure the flows' RTTs to derive the value for T_{max} in a similar way to measuring the mean service rate, but in the examples presented here we simply use a fixed value of 200ms since this is an approximate upper bound on the RTT of the majority of the current Internet flows.

We note that the classical BDP rule is derived from the behavior of TCP congestion control (in particular, the reduction of cwnd by half on packet loss) and assumes a constant service rate and fluid-like packet arrivals. Hence, for example, at low service rates the BDP rule suggests use of extremely small buffer sizes. However, in addition to accommodating TCP behavior, buffers have the additional role of absorbing short-term packet bursts and, in the case of wireless links, short-term fluctuations in packet service times. It is these latter effects that lead to the steep drop-off in throughput efficiency that can be observed in Fig. 3 when there are competing uploads (and so stochastic variations in packet service times due to channel contention, see Fig. 2.) plus small

Algorithm 1 Drop tail operation of the eBDP algorithm.

- 1: Set the target queuing delay T_{max} .
 - 2: Set the over-provision parameter c .
 - 3: **for** each incoming packet p **do**
 - 4: Calculate $Q_{eBDP} = \min(T_{max}/T_{serv} + c, Q_{max}^{eBDP})$ where T_{serv} is from MAC Algorithm 2.
 - 5: **if** current queue occupancy $< Q_{eBDP}$ **then**
 - 6: Put p into queue
 - 7: **else**
 - 8: Drop p .
 - 9: **end if**
 - 10: **end for**
-

buffer sizes. We therefore modify the eBDP update rule to $Q_{eBDP} = \min(T_{max}/T_{serv} + c, Q_{max}^{eBDP})$ where c is an over-provisioning amount to accommodate short-term fluctuations in service rate. Due to the complex nature of the service time process at a wireless station (which is coupled to the traffic arrivals etc at other stations in the WLAN) and of the TCP traffic arrival process (where feedback creates coupling to the service time process), obtaining an analytic value for c is intractable. Instead, based on the measurements in Fig. 3 and others, we have found empirically that a value of $c = 5$ packets works well across a wide range of network conditions. Pseudo-code for eBDP is shown in Algorithms 1 and 2.

The effectiveness of this simple adaptive algorithm is illustrated in Fig. 4. Fig. 4(a) shows the buffer size and queue occupancy time histories when only a single station is active in a WLAN while Fig. 4(b) shows the corresponding results when ten additional stations now also contend for channel access. Comparing with Fig. 3(e), it can be seen that buffer sizes of 330 packets and 70 packets, respectively, are needed to yield 100% throughput efficiency and eBDP selects buffer

Algorithm 2 MAC operation of the eBDP algorithm.

- 1: Set the averaging parameter W .
- 2: **for** each outgoing packet p **do**
- 3: Record service start time t_s for p .
- 4: Wait until receive MAC ACK for p , record service end time t_e .
- 5: Calculate service time of p : $T_{serv} = (1 - W)T_{serv} + W(t_e - t_s)$.
- 6: **end for**

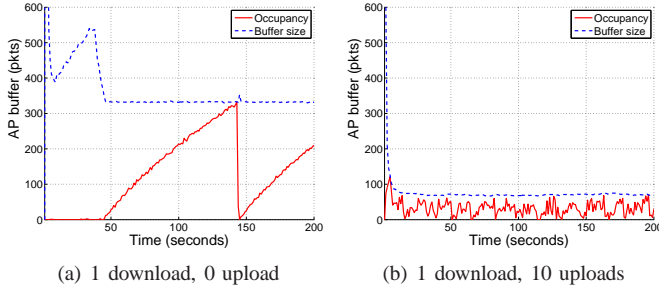


Fig. 4. Histories of buffer size and buffer occupancy with the eBDP algorithm. In (a) there is one download and no upload flows. In (b) there are 1 download and 10 upload flows. 54/6Mbps physical data/basic rates.

sizes which are in good agreement with these thresholds.

In Fig. 5 we plot the throughput efficiency (measured as the ratio of the achieved throughput to that with a fixed 400-packet buffer) and max smoothed RTT over a range of network conditions obtained using the eBDP algorithm. It can be seen that the adaptive algorithm maintains high throughput efficiency across the entire range of operating conditions. This is achieved while maintaining the latency approximately constant at around 400ms (200ms propagation delay plus $T_{max} = 200$ ms queuing delay) – the latency rises slightly with the number of uploads due to the over-provisioning parameter c used to accommodate stochastic fluctuations in service rate.

While $T_{max} = 200$ ms is used as the target drain time in the eBDP algorithm, realistic traffic tends to consist of flows with a mix of RTTs. Fig. 6 plots the results as we vary the RTT of the wired backhaul link while keeping $T_{max} = 200$ ms. We observe that the throughput efficiency is close to 100% for RTTs up to 200ms. For an RTT of 300ms, we observe a

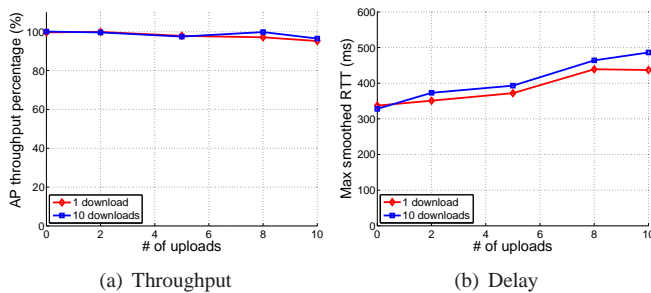


Fig. 5. Performance of the eBDP algorithm as the number of upload flows is varied. Data is shown for 1, 10 download flows and 0, 2, 5, 10 uploads. Wired RTT 200ms. Here the AP throughput percentage is the ratio between the throughput achieved using the eBDP algorithm and that by a fixed buffer size of 400 packets (i.e. the maximum achievable throughput in this case).

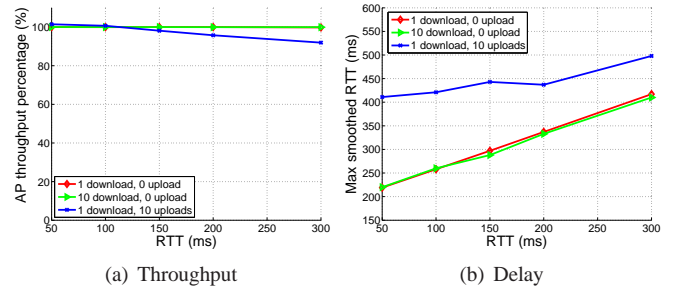


Fig. 6. Performance of the eBDP algorithm as the RTT of wired backhaul is varied. Data is shown for 1, 10 downloads and 0, 10 uploads. Here the AP throughput percentage is the ratio between the throughput achieved using the eBDP algorithm and that by a fixed buffer size of 400 packets (i.e. the maximum achievable throughput in this case).

slight decrease in throughput when there is 1 download and 10 contending upload flows, which is to be expected since T_{max} is less than the link delay and so the buffer is less than the BDP. This could be improved by measuring the average RTT instead of using a fixed value, but it is not clear that the benefit is worth the extra effort. We also observe that there is a difference between the max smoothed RTT with and without upload flows. The RTT in our setup consists of the wired link RTT, the queuing delays for TCP data and ACK packets and the MAC layer transmission delays for TCP data and ACK packets. When there are no upload flows, TCP ACK packets can be transmitted with negligible queuing delays since they only have to contend with the AP. When there are upload flows however, stations with TCP ACK packets have to contend with other stations sending TCP data packets as well. TCP ACK packets therefore can be delayed accordingly, which causes the increase in RTT observed in Fig. 6.

Fig. 7 demonstrates the ability of the eBDP algorithm to respond to changing network conditions. At time 300s the number of uploads is increased from 0 to 10 flows. It can be seen that the buffer size quickly adapts to the changed conditions when the weight $W = 0.001$. This roughly corresponds to averaging over the last 1000 packets². When the number of uploads is increased at time 300s, it takes 0.6 seconds (current throughput is 13.5Mbps so $t = 1000 * 8000 / 13.5 * 10^6 = 0.6$) to send 1000 packets, i.e., the eBDP algorithm is able to react to network changes roughly on a timescale of 0.6 second.

V. EXPLOITING STATISTICAL MULTIPLEXING: THE A* ALGORITHM

While the eBDP algorithm is simple and effective, it is unable to take advantage of the statistical multiplexing of TCP cwnd backoffs when multiple flows share the same link. For example, it can be seen from Fig. 8 that while a buffer size of 338 packets is needed to maximize throughput with a single download flow, this falls to around 100 packets when 10 download flows share the link. However, in both cases the eBDP algorithm selects a buffer size of approximately 350

²As per [8], the current value is averaged over the last t observations for $x\%$ percentage of accuracy where $x = 1 - (1 - W)^t$, t is the number of updates (which are packets in our case). When $W = 0.001$ and $t = 1000$ we have that $x = 0.64$.

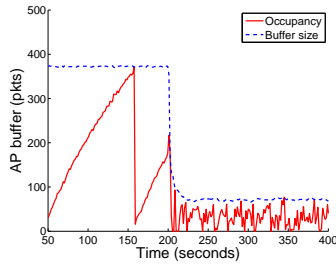


Fig. 7. Convergence of the eBDP algorithm following a change in network conditions. One download flow. At time 200s the number of upload flows is increased from 0 to 10.

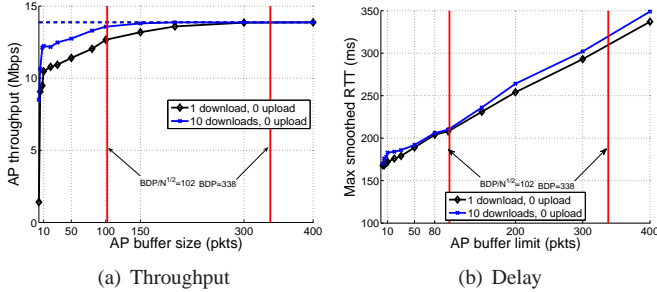


Fig. 8. Impact of statistical multiplexing. There are 1/10 downloads and no uploads. Wired RTT 200ms.

packets (see Figs. 4(a) and 9). It can be seen from Fig. 9 that as a result with the eBDP algorithm the buffer rarely empties when 10 flows share the link. That is, the potential exists to lower the buffer size without loss of throughput.

In this section we consider the design of a measurement-based algorithm (the ALT algorithm) that is capable of taking advantage of such statistical multiplexing opportunities.

A. Adaptive Limit Tuning (ALT) Feedback Algorithm

Our objective is to simultaneously achieve both efficient link utilization and low delays in the face of stochastic time-variations in the service time. Intuitively, for efficient link utilization we need to ensure that there is a packet available to transmit whenever the station wins a transmission opportunity. That is, we want to minimize the time that the station buffer lies empty, which in turn can be achieved by making the buffer size sufficiently large (under fairly general traffic conditions, buffer occupancy is a monotonically increasing function of

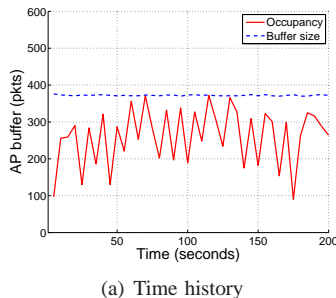


Fig. 9. Histories of buffer size and buffer occupancy with the eBDP algorithm when there are 10 downloads and no uploads.

buffer size [13].). However, using large buffers can lead to high queuing delays, and to ensure low delays the buffer should be as small as possible. We would therefore like to operate with the smallest buffer size that ensures sufficiently high link utilization. This intuition suggests the following approach. We observe the buffer occupancy over an interval of time. If the buffer rarely empties, we decrease the buffer size. Conversely, if the buffer is observed to be empty for too long, we increase the buffer size. Of course, further work is required to convert this basic intuition into a well-behaved algorithm suited to practical implementation. Not only do the terms “rarely”, “too long” etc need to be made precise, but we note that an inner feedback loop is created whereby buffer size is adjusted depending on the measured link utilization, which in turn depends on the buffer size. This new feedback loop is in addition to the existing outer feedback loop created by TCP congestion control, whereby the offered load is adjusted based on the packet loss rate, which in turn is dependent on buffer size. Stability analysis of these cascaded loops is therefore essential.

We now introduce the following Adaptive Limit Tuning (ALT) algorithm. The dynamics and stability of this algorithm will then be analyzed in later sections. Define a queue occupancy threshold q_{thr} and let $t_i(k)$ (referred to as the *idle time*) be the duration of time that the queue spends at or below this threshold in a fixed observation interval t , and $t_b(k)$ (referred to as the *busy time*) be the corresponding duration spent above the threshold. Note that $t = t_i(k) + t_b(k)$ and the aggregate amount of idle/busy time t_i and t_b over an interval can be readily observed by a station. Also, the link utilisation is lower bounded by $t_b/(t_b + t_i)$. Let $q(k)$ denote the buffer size during the k -th observation interval. The buffer size is then updated according to

$$q(k+1) = q(k) + a_1 t_i(k) - b_1 t_b(k), \quad (1)$$

where a_1 and b_1 are design parameters. Pseudo-code for this ALT algorithm is given in Algorithm 3. This algorithm seeks to maintain a balance between the time t_i that the queue is idle and the time t_b that the queue is busy. That is, it can be seen that when $a_1 t_i(k) = b_1 t_b(k)$, the buffer size is kept unchanged. When the idle time is larger so that $a_1 t_i(k) > b_1 t_b(k)$, then the buffer size is increased. Conversely, when the busy time is large enough that $a_1 t_i(k) < b_1 t_b(k)$, then the buffer size is decreased.

More generally, assuming q converges to a stationary distribution (we discuss this in more detail later), then in steady-state we have that $a_1 E[t_i] = b_1 E[t_b]$, i.e., $E[t_i] = \frac{b_1}{a_1} E[t_b]$ and the mean link utilization is therefore lower bounded by

$$E\left[\frac{t_b}{t_i + t_b}\right] = \frac{E[t_b]}{t} = \frac{1}{1 + b_1/a_1}. \quad (2)$$

where we have made use of the fact that $t = t_i(k) + t_b(k)$ is constant. It can therefore be seen that choosing $\frac{b_1}{a_1}$ to be small then ensures high utilization. Choosing values for the parameters a_1 and b_1 is discussed in detail in Section V-B, but we note here that values of $a_1 = 10$ and $b_1 = 1$ are found to work well and unless otherwise stated are used in this paper. With regard to the choice of observation interval

Algorithm 3 : The ALT algorithm.

- 1: Set the initial queue size, the maximum buffer size q_{max} and the minimum buffer size q_{min} .
- 2: Set the increase step size a_1 and the decrease step size b_1 .
- 3: **for** Every t seconds **do**
- 4: Measure the idle time t_i .
- 5: $q_{ALT} = q_{ALT} + a_1 t_i - b_1 (t - t_i)$.
- 6: $q_{ALT} = \min(\max(q_{ALT}, q_{min}), q_{max})$
- 7: **end for**

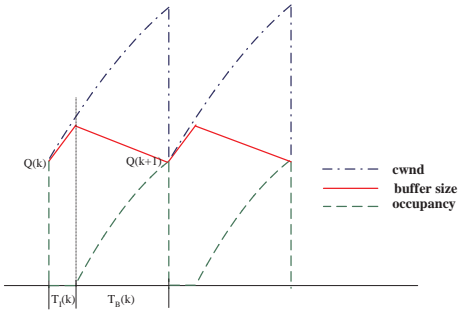


Fig. 10. Illustrating evolution of the buffer size.

t , this is largely determined by the time required to obtain accurate estimates of the queue idle and busy times. In the remainder of this paper we find a value of $t = 1$ second to be a good choice.

It is prudent to constrain the buffer size q to lie between the minimum and the maximum values q_{min} and q_{max} . In the following, the maximum size q_{max} and the minimum buffer size q_{min} are set to be 1600 and 5 packets respectively.

B. Selecting the Step Sizes for ALT

Define a *congestion event* as an event where the sum of all senders' TCP cwnd decreases. This cwnd decrease can be caused by the response of TCP congestion control to a single packet loss, or multiple packet losses that are lumped together in one RTT. Define a *congestion epoch* as the duration between two adjacent congestion events.

Let $Q(k)$ denote the buffer size at the k -th congestion event. Then,

$$Q(k+1) = Q(k) + aT_I(k) - bT_B(k) \quad (3)$$

where T_I is the "idle" time, i.e., the duration in seconds when the queue occupancy is below q_{thr} during the k -th congestion epoch, and T_B the "busy" time, i.e., the duration when the queue occupancy is above q_{thr} . This is illustrated in Fig. 10 for the case of a single TCP flow.

Notice that $a = a_1$ and $b = b_1$ where a_1 and b_1 are parameters used in the ALT algorithm. In the remainder of this section we investigate conditions to guarantee convergence and stability of the buffer dynamics with TCP traffic, which naturally lead to guidelines for the selection of a_1 and b_1 . We first define some TCP related quantities before proceeding.

Consider the case where TCP flows may have different round-trip times and drops need not be synchronized. Let n

be the number of TCP flows sharing a link, $w_i(k)$ be the cwnd of flow i at the k -th congestion event, T_i the round-trip propagation delay of flow i . To describe the cwnd additive increase we define the following quantities: (i) α_i is the rate in *packet/s* at which flow i increases its congestion window³, (ii) $\alpha_T = \sum_{i=1}^n \alpha_i$ is the aggregate rate at which flows increase their congestion windows, in packets/s, and (iii) $A_T = \sum_{i=1}^n \alpha_i / T_i$ approximates the aggregate rate, in *packets/s²*, at which flows increase their sending rate. Following the k -th congestion event, flows backoff their cwnd to $\beta_i(k)w_i(k)$. Flows may be unsynchronized, i.e., not all flows need back off at a congestion event. We capture this with $\beta_i(k) = 1$ if flow i does not backoff at event k . We assume that the α_i are constant and that the $\beta_i(k)$ (i.e. the pattern of flow backoffs) are independent of the flow congestion windows $w_i(k)$ and the buffer size $Q(k)$ (this appears to be a good approximation in many practical situations, see [26]).

To relate the queue occupancy to the flow cwnds, we adopt a fluid-like approach and ignore sub-RTT burstiness. We also assume that q_{thr} is sufficiently small relative to the buffer size that we can approximate it as zero. Considering now the idle time $T_I(k)$, on backoff after the k -th congestion event, if the queue occupancy does not fall below q_{thr} then $T_I(k) = 0$. Otherwise, immediately after backoff the send rate of flow i is $\beta_i(k)w_i(k)/T_i$ and we have that

$$T_I(k) = \frac{E[B] - \sum_{i=1}^n \beta_i(k)w_i(k)/T_i}{A_T}, \quad (4)$$

where $E[B]$ is the mean service rate of the considered buffer.

At congestion event k the aggregate flow throughput necessarily equals the link capacity, i.e.,

$$\sum_{i=1}^n \frac{w_i(k)}{T_i + Q(k)/E[B]} = E[B].$$

We then have that

$$\begin{aligned} \sum_{i=1}^n \frac{w_i(k)}{T_i} &= \sum_{i=1}^n \frac{w_i(k) T_i + Q(k)/E[B]}{T_i T_i + Q(k)/E[B]} \\ &= \sum_{i=1}^n \frac{w_i(k)}{T_i + Q(k)/E[B]} + \\ &\quad \frac{Q(k)}{E[B]} \sum_{i=1}^n \frac{w_i(k)}{T_i + Q(k)/E[B]} \frac{1}{T_i} \end{aligned}$$

Assume that the spread in flow round-trip propagation delays and congestion windows is small enough that $\sum_{i=1}^n (w_i(k)/(E[B]T_i + Q(k))(1/T_i))$ can be accurately approximated by $1/T_T$, where $T_T = \frac{n}{\sum_{i=1}^n \frac{1}{T_i}}$ is the harmonic mean of T_i . Then

$$\sum_{i=1}^n \frac{w_i(k)}{T_i} \approx E[B] + \frac{Q(k)}{T_T},$$

and

$$T_I(k) \approx \frac{(1 - \beta_T(k))E[B] - \beta_T(k)Q(k)/T_T}{A_T} \quad (5)$$

³Standard TCP increases the flow congestion window by one packet per RTT, in which case $\alpha_i \approx 1/T_i$.

where $\beta_T(k) = \frac{\sum_{i=1}^n \beta_i(k) w_i(k) / T_i}{\sum_{i=1}^n w_i(k) / T_i}$ is the effective aggregate backoff factor of the flows. When flows are synchronized, i.e., $\beta_i = \beta \forall i$, then $\beta_T = \beta$. When flows are unsynchronized but have the same *average* backoff factor, i.e., $E[\beta_i] = \beta$, then $E[\beta_T] = \beta$.

If the queue empties after backoff, the queue busy time $T_B(k)$ is directly given by

$$T_B(k) = Q(k+1)/\alpha_T \quad (6)$$

where $\alpha_T = \sum_{i=1}^n \alpha_i$ is the aggregate rate at which flows increase their congestion windows, in packets/s. Otherwise,

$$T_B(k) = (Q(k+1) - q(k))/\alpha_T \quad (7)$$

where $q(k)$ is the buffer occupancy after backoff. It turns out that for the analysis of stability it is not necessary to calculate $q(k)$ explicitly. Instead, letting $\delta(k) = q(k)/Q(k)$, it is enough to note that $0 \leq \delta(k) < 1$.

Combining (3), (5), (6) and (7),

$$Q(k+1) = \begin{cases} \lambda_e(k)Q(k) + \gamma_e(k)E[B]T_T, & q(k) \leq q_{thr} \\ \lambda_f(k)Q(k), & \text{otherwise} \end{cases}$$

where

$$\begin{aligned} \lambda_e(k) &= \frac{\alpha_T - a\beta_T(k)\alpha_T/(A_T T_T)}{\alpha_T + b}, \\ \lambda_f(k) &= \frac{\alpha_T + b\delta(k)}{\alpha_T + b}, \quad \gamma_e(k) = a \frac{1 - \beta_T(k)}{\alpha_T + b} \frac{\alpha_T}{A_T T_T}. \end{aligned}$$

Taking expectations,

$$\begin{aligned} E[Q(k+1)] &= E[\lambda_e(k)Q(k) + \gamma_e(k)E[B]T_T | q(k) \leq q_{thr}] p_e(k) \\ &\quad + E[\lambda_f(k)Q(k) | q(k) > q_{thr}] (1 - p_e(k)) \end{aligned}$$

with $p_e(k)$ the probability that the queue $\leq q_{thr}$ following the k -th congestion event. Since the $\beta_i(k)$ are assumed independent of $Q(k)$ we may assume that $E[Q(k) | q(k) \leq q_{thr}] = E[Q(k) | q(k) > q_{thr}] = E[Q(k)]$ and

$$E[Q(k+1)] = \lambda(k)E[Q(k)] + \gamma(k)E[B]T_T \quad (8)$$

where

$$\begin{aligned} \lambda(k) &= p_e(k)E[\lambda_e(k) | q(k) \leq q_{thr}] \\ &\quad + (1 - p_e(k))E[\lambda_f(k) | q(k) > q_{thr}], \\ \gamma(k) &= p_e(k)E[\gamma_e(k) | q(k) \leq q_{thr}] \end{aligned}$$

C. A Sufficient Condition for Stability

Provided $|\lambda(k)| < 1$ the queue dynamics in (8) are exponentially stable. In more detail, $\lambda(k)$ is the convex combination of $E[\lambda_e(k)]$ and $E[\lambda_f(k)]$ (where the conditional dependence of these expectations is understood, but omitted to streamline notation). Stability is therefore guaranteed provided $|E[\lambda_e(k)]| < 1$ and $|E[\lambda_f(k)]| < 1$. We have that $0 < E[\lambda_f(k)] < 1$ when $b > 0$ since α_T is non-negative and $0 \leq \delta(k) < 1$. The stability condition is therefore that $|E[\lambda_e(k)]| < 1$.

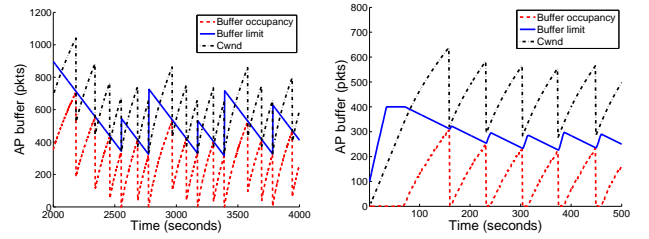


Fig. 11. Instability and stability of the ALT algorithm. In (a), $a=100$, $b=1$, the maximum buffer size is 50000 packets. In (b), $a=10$, $b=1$, the maximum buffer size is 400 packets. In both figures, there is 1 download and no upload.

Under mild independence conditions,

$$E[\lambda_e(k)] = \frac{\alpha_T - aE[\beta_T(k)]\alpha_T/(A_T T_T)}{\alpha_T + b}.$$

Observe that,

$$\frac{\alpha_T}{A_T T_T} = \frac{1}{n} \frac{(\sum_{i=1}^n 1/T_i)^2}{\sum_{i=1}^n 1/T_i^2}$$

when we use the standard TCP AIMD increase of one packet per RTT, in which case $\alpha_i \approx 1/T_i$. We therefore have that $1/n \leq \alpha_T/(A_T T_T) \leq 1$. Also, when the standard AIMD backoff factor of 0.5 is used, $0.5 < E[\beta_T(k)] < 1$. Thus, since $a > 0$, $b > 0$, $\alpha_T > 0$, it is sufficient that

$$-1 < \frac{\alpha_T - a}{\alpha_T + b} \leq E[\lambda_e(k)] \leq \frac{\alpha_T}{\alpha_T + b} < 1$$

A sufficient condition (from the left inequality) for stability is then that $a < 2\alpha_T + b$. Using again (as in the eBDP algorithm) 200ms as the maximum RTT, a rough lower bound on α_T is 5 (corresponding to 1 flow with RTT 200ms). The stability constraint is then that

$$a < 10 + b. \quad (9)$$

Fig. 11(a) demonstrates that the instability is indeed observed in simulations. Here, $a = 100$ and $b = 1$ are used as example values, i.e., the stability conditions are not satisfied. It can be seen that the buffer size at congestion events oscillates around 400 packets rather than converging to a constant value. We note, however, that in this example and others the instability consistently manifests itself in a benign manner (small oscillations). However, we leave detailed analysis of the onset of instability as future work.

Fig. 11(b) shows the corresponding results with $a = 10$ and $b = 1$, i.e., when the stability conditions are satisfied. It can be seen that the buffer size at congestion events settles to a constant value, thus the buffer size time history converges to a periodic cycle.

D. Fixed point

When the system dynamics are stable and $p_e = 0$, from (8) we have that

$$\lim_{k \rightarrow \infty} E[Q(k)] = \frac{(1 - E[\beta_T])}{b/a + E[\beta_T]} E[B]T_T. \quad (10)$$

For synchronized flows with the standard TCP backoff factor of 0.5 (i.e., $E[\beta_T] = 0.5$) and the same RTT,

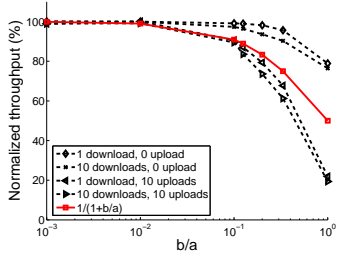


Fig. 12. Impact of b/a on throughput efficiency. The maximum buffer size is 400 packets, and the minimum buffer size is 2 packets.

$\frac{(1-E[\beta_T])}{b/a+E[\beta_T]}E[B]T_T$ reduces to the BDP when $b/a = 0$. This indicates that for high link utilization we would like the ratio b/a to be small. Using (5), (6) and (10) we have that in steady-state the expected link utilization is lower bounded by

$$\frac{1}{1 + \frac{b}{a} \frac{\alpha_T}{A_T T_T}} \geq \frac{1}{1 + \frac{b}{a}}. \quad (11)$$

This lower bound is plotted in Fig. 12 together with the measured throughput efficiency vs b/a in a variety of traffic conditions. Note that in this figure the lower bound is violated by the measured data when $b/a > 0.1$ and we have a large number of uploads. At such large values of b/a plus many contending stations, the target buffer sizes are extremely small and micro-scale burstiness means that TCP RTOs occur frequently. It is this that leads to violation of the lower bound (11) (since $p_e = 0$ does not hold). However, this corresponds to an extreme operating regime and for smaller values of b/a the lower bound is respected. It can be seen from Fig. 12 that the efficiency decreases when the ratio of b/a increases. In order to ensure throughput efficiency $\geq 90\%$ it is required that

$$\frac{b}{a} \leq 0.1. \quad (12)$$

Combined with the stability condition in inequality (9), we have that $a = 10, b = 1$ are feasible integer values, that is, we choose $a_1 = 10$ and $b_1 = 1$ for the A* algorithm.

E. Convergence rate

In Fig. 13(a) we illustrate the convergence rate of the ALT algorithm. There is one download, and at time 500s the number of upload flows is increased from 0 to 10. It can be seen that the buffer size limit converges to its new value in around 200 seconds or 3 minutes. In general, the convergence rate is determined by the product $\lambda(0)\lambda(1)\dots\lambda(k)$. In this example, the buffer does not empty after backoff and the convergence rate is thus determined by $\lambda_f(k) = \frac{\alpha_T + b\delta(k)}{\alpha_T + b}$. To achieve fast convergence, we require small $\lambda_f(k)$ so that $Q(k+1) = \lambda_f(k)Q(k)$ is decreased quickly to the desired value. We thus need large b to achieve fast convergence. However, $b = 1$ is used here in order to respect the stability condition in (9) and the throughput efficiency condition in (12). Note that when conditions change such that the buffer size needs to increase, the convergence rate is determined by the a parameter. This has a value of $a = 10$ and thus the algorithm adapts much more quickly to increase the buffer than to decrease it and the

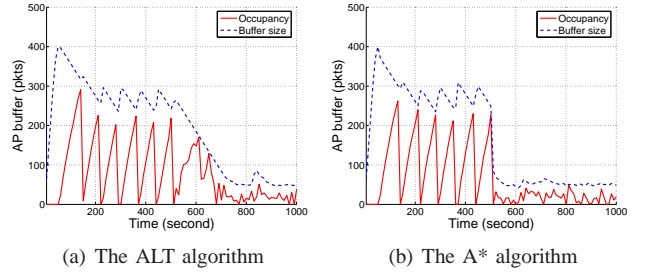


Fig. 13. Convergence rate of the ALT and A* algorithms. One download flow, $a = 10, b = 1$. At time 500s the number of upload flows is increased from 0 to 10.

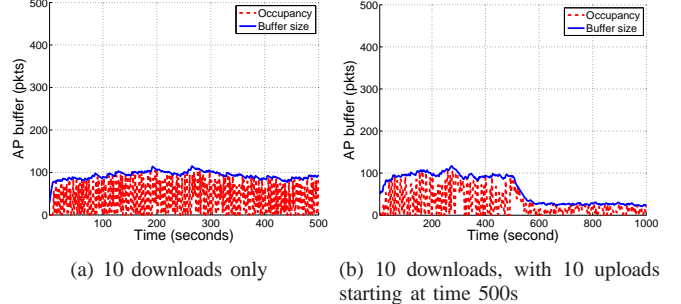


Fig. 14. Buffer time histories with the A* algorithm, $a=10, b=1$.

example in Fig. 13(a) is essentially a worst case. In the next section, we address the slow convergence by combining the ALT and the eBDP algorithms to create a hybrid algorithm.

F. Combining eBDP and ALT: The A* Algorithm

We can combine the eBDP and ALT algorithms by using the mean packet service time to calculate Q_{eBDP} as per the eBDP algorithm (see Section IV), and the idle/busy times to calculate q_{ALT} as per the ALT algorithm. We then select the buffer size as $\min\{Q_{eBDP}, q_{ALT}\}$ to yield a hybrid algorithm, referred to as the A* algorithm, that combines the eBDP and the ALT algorithms.

When channel conditions change, the A* algorithm uses the eBDP measured service time to adjust the buffer size promptly. The convergence rate depends on the smoothing weight W . As calculated in Section IV, it takes around 0.6 second for Q_{eBDP} to converge. The A* algorithm can further use the ALT algorithm to fine tune the buffer size to exploit the potential reduction due to statistical multiplexing. The effectiveness of this hybrid approach when the traffic load is increased suddenly is illustrated in Fig. 13(b) (which can be directly compared with Fig. 13(a)). Fig. 14(b) shows the corresponding time histories for 10 download flows and a changing number of competing uploads.

G. Performance

The basic impetus for the design of the A* algorithm is to exploit the possibility of statistical multiplexing to reduce buffer sizes. Fig. 14(a) illustrates the performance of the A* algorithm when there are 10 downloads and no upload flows. Comparing with the results in Fig. 9 using fixed size buffers, we can see that the A* algorithm can achieve significantly

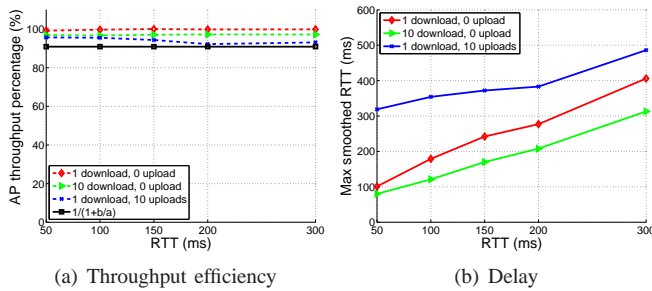


Fig. 16. Performance of the A* algorithm as the wired RTT is varied. Physical layer data and basic rates are 54 and 6 Mbps. Here the AP throughput percentage is the ratio between the throughput achieved using the A* algorithm and the maximum throughput using fixed size buffers.

smaller buffer sizes (i.e., a reduction from more than 350 packets to 100 packets approximately) when multiplexing exists. Fig. 15 summarizes the throughput and delay performance of the A* algorithm for a range of network conditions (numbers of uploads and downloads) and physical transmit rates ranging from 1Mbps to 216Mbps. This can be compared with Fig. 3. It can be seen that in comparison with the use of a fixed buffer size the A* algorithm is able to achieve high throughput efficiency across a wide range of operating conditions while minimizing queuing delays.

In Fig. 16 we further evaluate the A* algorithm when the wired RTTs are varied from 50-300ms and the number of uploads is varied from 0-10. Comparing these with the results (Fig. 5 and 6) of the eBDP algorithm we can see that the A* algorithm is capable of exploiting the statistical multiplexing where feasible. In particular, significantly lower delays are achieved with 10 download flows whilst maintaining comparable throughput efficiency.

H. Impact of Channel Errors

In the foregoing simulations the channel is error free and packet losses are solely due to buffer overflow and MAC-layer collisions. In fact, channel errors have only a minor impact on the effectiveness of buffer sizing algorithms as errors play a similar role to collisions with regard to their impact on link utilization. We support this claim first using a simulation example with a channel having an i.i.d noise inducing a bit error rate (BER) of 10^{-5} . Results are shown in Fig. 17 where we can see a similar trend as in the cases when the medium is error free (Figs. 15(e) 15(f)).

We further confirm this claim in our test-bed implementations where tests were conducted in 802.11b/g channels and noise related losses were observed. See Section VI for details.

I. DCF Operation

The proposed buffer sizing algorithms are still valid for DCF since link utilization and delay considerations remain applicable, as is the availability of service time (for the eBDP algorithm) and idle/busy time measurements (for the ALT algorithm). In particular, if the considered buffer is heavily backlogged, to ensure low delays, the buffer size should be reduced. If otherwise the buffer lies empty, it may be due

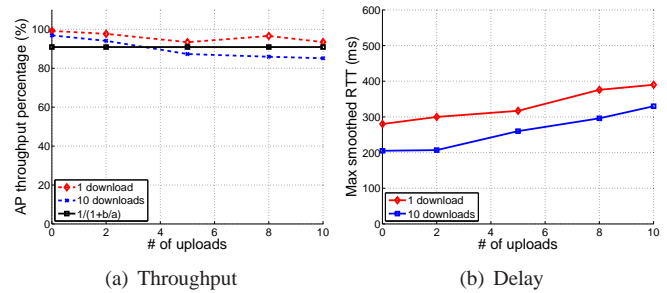


Fig. 17. Performance of the A* algorithm when the channel has a BER of 10^{-5} . Physical layer data and basic rates are 54 and 6 Mbps. Here the AP throughput percentage is the ratio between the throughput achieved using the A* algorithm and the maximum throughput using fixed size buffers.

to that the current buffer size is too small which causes the TCP source back off after buffer overflow. To accommodate more future packets, the buffer size can be increased. Note that increasing buffer sizes in this case would not lead to high delays but has the potential to improve throughput. This tradeoff between the throughput and the delays thus holds for both EDCA and DCF.

However, the DCF allocates roughly equal numbers of transmission opportunities to stations. A consequence of using DCF is thus that when the number of upload flows increases, the uploads may produce enough TCP ACK packets to keep the AP's queue saturated. In fact, once there are two upload flows, TCP becomes unstable due to repeated timeouts (see [20] for a detailed demonstration), causing the unfairness issue discussed in Section II-C. Therefore, we present results for up to two uploads in Fig. 18, as this is the greatest number of upload flows where TCP with DCF can exhibit stable behavior using both fixed size buffers and the A* algorithm. Note that in this case using the A* algorithm on upload stations can also decrease the delays and maintain high throughput efficiency if their buffers are frequently backlogged.

We also present results when there are download flows only (so the unfairness issue does not exist). Fig. 19 illustrates the throughput and delay performance achieved using the A* algorithm and fixed 400-packet buffers. As in the EDCA cases, we can see that the A* algorithm is able to maintain a high throughput efficiency with comparatively low delays.

Note that DCF is also used in the production WLAN test where the A* algorithm is observed to perform well (see Section I).

J. Rate Adaptation

We did not implement rate adaptation in our simulations. However, we did implement the A* algorithm in the Linux MadWifi driver which includes rate adaptation algorithms. We tested the A* algorithm in the production WLAN of the Hamilton Institute with the default SampleRate algorithm enabled. See Section I.

VI. EXPERIMENTAL RESULTS

We have implemented the proposed algorithms in the Linux MadWifi driver, and in this section we present tests on an

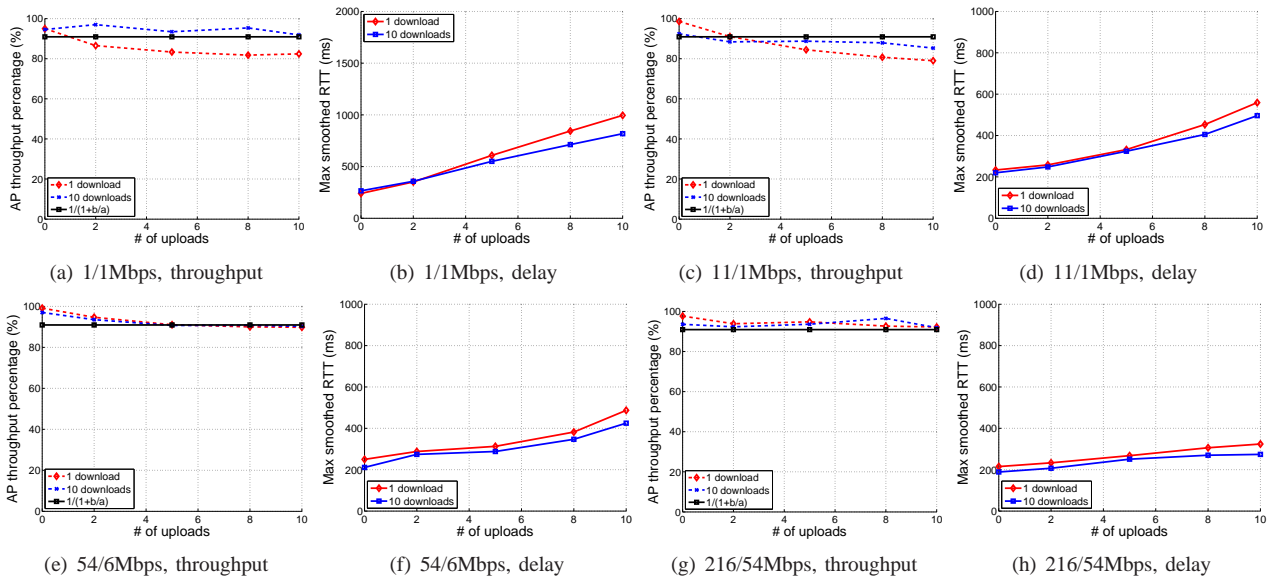


Fig. 15. Throughput efficiency and maximum smoothed round trip delays (max sRTT) for the topology in Fig. 1 when the A* algorithm is used. Here, the AP throughput efficiency is the ratio between the throughput achieved using the A* algorithm and the maximum throughput achieved using fixed size buffers. Rates before and after the '/' are used physical layer data and basic rates. For the 216Mbps data, 8 packets are aggregated into each frame at the MAC layer to improve throughput efficiency in an 802.11n-like scheme. The wired RTT is 200 ms.

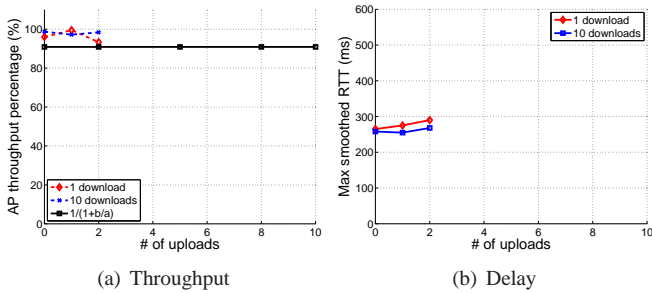


Fig. 18. Performance of the A* algorithm for 802.11 DCF operation when there are both upload and download flows in the network. Here the AP throughput percentage is the ratio between the throughput achieved using the A* algorithm and the maximum throughput using fixed size buffers. Physical layer data and basic rates used are 54 and 6 Mbps.

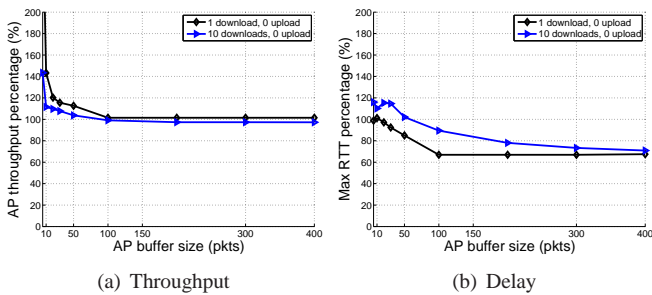


Fig. 19. Performance of the A* algorithm for 802.11 DCF operation when there are download flows only in the network. Here we illustrate the percentage (between the results achieved using the A* algorithm and those using varied AP buffer sizes as shown on the x-axis) of both throughput and delays. Physical layer data and basic rates used are 54 and 6 Mbps.

experimental testbed located in an office environment and introduce results illustrating operation with complex traffic that includes both TCP and UDP, a mix of uploads and downloads, and a mix of connection sizes.

A. Testbed Experiment

The testbed topology is shown in Fig. 20. A wired network is emulated using a desktop PC running dummynet software on FreeBSD 6.2 which enables link rates and propagation delays to be controlled. The wireless AP and the server are connected to the dummynet PC by 100Mbps Ethernet links. Routing in the network is statically configured. Network management is carried out using ssh over a wired control plane to avoid affecting wireless traffic.

In the WLAN, a desktop PC is used as the AP and 12 PC-based embedded Linux boxes based on the Soekris net4801 are used as client stations. All are equipped with an Atheros 802.11b/g PCI card with an external antenna. All nodes run a Linux 2.6.21.1 kernel and a MadWifi wireless driver (version r2366) modified to allow us to adjust the 802.11e CW_{min} , CW_{max} and $AIFS$ parameters as required. Specific vendor features on the wireless card, such as turbo mode, rate adaptation and multi-rate retries, are disabled. All of the tests are performed using a transmission rate of 11Mbps (i.e., we use an 802.11b PHY) with RTS/CTS disabled and the channel number explicitly set. Channel 1 has been selected to carry out the experiments. The testbed is not in an isolated radio environment, and is subject to the usual impairments seen in an office environment. Since the wireless stations are based on low power embedded systems, we have tested these wireless stations to confirm that the hardware performance (especially the CPU) is not a bottleneck for wireless transmissions at the 11Mbps PHY rate used. The configuration of the various network buffers and MAC parameters is detailed in Table II.

Although both SACK enabled TCP NewReno and TCP CUBIC with receiver buffers of 4096KB have been tested, here we only report the results for the latter as CUBIC is now the default congestion control algorithm used in Linux. Default values of Linux Kernel 2.6.21.1 are used for all the

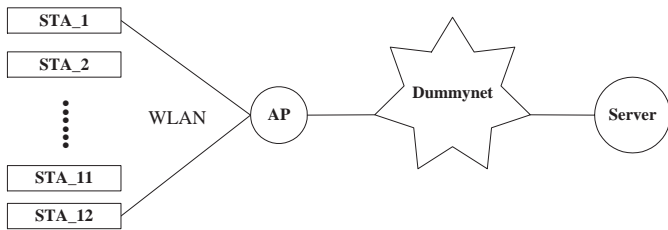


Fig. 20. Topology used in experimental tests.

Parameters	Values
Interface tx queue	2 packets
Dummynet queue	100 packets
MAC Preamble	long
MAC data rate	11Mbps
MAC ACK rate	11Mbps
MAC retries	11

TABLE II
TESTBED PARAMETERS SUMMARY.

other TCP parameters. We put TCP ACK packets into a high priority queue (we use the WME_AC_VO queue of MadWifi as an example) which is assigned parameters of $CW_{min} = 3$, $CW_{max} = 7$ and $AIFS = 2$. TCP data packets are collected into a lower priority queue (we use the WME_AC_VI queue) which is assigned $CW_{min} = 31$, $CW_{max} = 1023$ and $AIFS = 6$. We use iperf to generate TCP traffic and results are collected using both iperf and tcpdump.

B. Traffic Mix

We configure the traffic mix on the network to capture the complexity of real networks in order to help gain greater confidence in the practical utility of the proposed buffer sizing approach. With reference to Fig. 20 showing the network topology, we create the following traffic flows:

- *TCP uploads.* One long-lived TCP upload from each of STAs 1, 2 and 3 to the server in the wired network. STAs 2 and 3 always use a fixed 400-packet buffer, while STA 1 uses both a fixed 400-packet buffer and the A* algorithm.
- *TCP downloads.* One long-lived TCP download from the wired server to each of STAs 4, 5 and 6.
- *Two way UDP.* One two-way UDP flow from the wired server to STA 7. The packet size used is 64 bytes and the mean inter-packet interval is 1s. Another UDP flow from the wired server to STA 8 with the used packet size of 1000 bytes and the mean inter-packet interval of 1s.
- *Mix of TCP connection sizes.* These flows mimic web traffic⁴. A short TCP download from the wired server to STA 9, the connection size of which is 5KB (approximately 3 packets). A slightly longer TCP download from the wired server to STA 10 with a connection size of 20KB (approximately 13 packets) and another to STA 11 (connection size 30KB, namely, around 20 packets). A fourth connection sized 100KB from the server to STA 12. For each size of these connection, a new flow is

started every 10s to allow collection of statistics on the mean completion time.

C. Results

Fig. 21 shows example time histories of the buffer size and occupancy at the AP with a fixed buffer size of 400 packets and when the A* algorithm is used for dynamic buffer sizing. Note that in this example the 400 packet buffer never completely fills. Instead the buffer occupancy has a peak value of around 250 packets. This is due to non-congestive packet losses caused by channel noise (the testbed operates in a real office environment with significant interference, because there are bluetooth devices and WLANs working in channel 1.) which prevent the TCP congestion window from growing to completely fill the buffer. Nevertheless, it can be seen that the buffer rarely empties and thus it is sufficient to provide an indication of the throughput when the wireless link is fully utilized.

We observe that while buffer histories are very different with a fixed size buffer and the A* algorithm, the throughput is very similar in these two cases (see Table III).

One immediate benefit of using smaller buffers is thus a reduction in network delays. Table IV shows the measured delays experienced by the UDP flows sharing the WLAN with the TCP traffic. It can be seen that for STA 8 both the mean and the maximum delays are significantly reduced when the A* algorithm is used. This potentially has major implications for time sensitive traffic when sharing a wireless link with data traffic. Note that the queuing delays from STA 7 are for traffic passing through the high-priority traffic class used for TCP ACKs, while the measurements from STA 8 are for traffic in the same class as TCP data packets. For the offered loads used, the service rate of the high-priority class is sufficient to avoid queue buildup and this is reflected in the measurements.

The reduction in network delay not only benefits UDP traffic, but also short-lived TCP connections. Fig. 22 shows the measured completion time vs connection size for TCP flows. It can be seen that the completion time is consistently lower by a factor of at least two when A* dynamic buffer sizing is used. Since the majority of internet flows are short-lived TCP connections (e.g., most web traffic), this potentially translates into a significant improvement in user experience.

Note that STA's 2 and 3 in the A* column of Table III use fixed size buffers rather than the A* algorithm. The results shown are the throughput they achieve when other stations run the A* algorithm. It can be seen that the A* algorithm does not significantly impact STAs 2 and 3, confirming that A* can support incremental roll-out without negatively impacting legacy stations that are using fixed size buffers.

VII. RELATED WORK

The classical approach to sizing Internet router buffers is the BDP rule proposed in [31]. Recently, in [5] it is argued that the BDP rule may be overly conservative on links shared by a large number of flows. In this case it is unlikely that TCP congestion window sizes (cwnd) evolve synchronously and due to statistical multiplexing of cwnd backoff, the combined

⁴Note that in the production WLAN test, we used real web traffic.

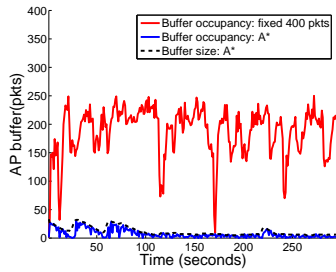


Fig. 21. Buffer size and occupancy time histories measured at the AP with fixed 400-packet buffers and the A* algorithm.

	Fixed 400 packets	A*
Throughput of STA 1	1.36Mbps	1.33Mbps
Throughput of STA 2	1.29Mbps	1.30Mbps
Throughput of STA 3	1.37Mbps	1.33Mbps
Throughput of STA 4	0.35Mbps	0.41Mbps
Throughput of STA 5	0.39Mbps	0.39Mbps
Throughput of STA 6	0.52Mbps	0.42Mbps

TABLE III
MEASURED THROUGHPUT.

buffer requirement can be considerably less than the BDP. The analysis in [5] suggests that it may be sufficient to size buffers as BDP/\sqrt{n} . This work is extended in [25], [10] and [33] to consider the performance of TCP congestion control with many connections under the assumption of small, medium and large buffer sizes. Several authors have pointed out that the value n can be difficult to determine for realistic traffic patterns, which not only include a mix of connections sizes and RTTs, but can also be strongly time-varying [9], [32]. In [32], it is observed from measurements on a production link that traffic patterns vary significantly over time, and may contain a complex mix of flow connection lengths and RTTs. It is demonstrated in [9][32] that the use of very small buffers can lead to an excessive loss rate. Motivated by these observations,

	Fixed 400 packets	A*
	mean (max)	mean (max)
RTT to STA 7	201ms (239ms)	200ms (236ms)
RTT to STA 8	1465ms (2430ms)	258ms (482ms)

TABLE IV
MEASURED DELAYS OF THE UDP FLOWS. STA 7’S TRAFFIC IS PRIORITIZED TO AVOID QUEUE BUILDUP AND THIS IS REFLECTED IN THE MEASUREMENTS.

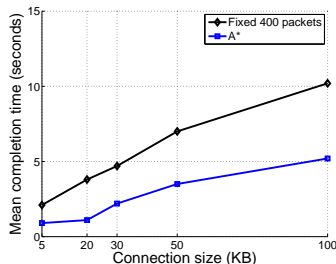


Fig. 22. Measured completion time vs connection size. Results are averages of multiple runs.

in [27] [12] a measurement-based adaptive buffer size tuning method is proposed. However, this approach is not applicable to WLANs since it requires a priori knowledge of the link capacity or line rate, which in WLANs is time-varying and load dependent. [34] introduces another adaptive buffer sizing algorithm based on control theory for Internet core routers. [24], [14] consider the role of the output/input capacity ratio at a network link in determining the required buffer size. [6] experimentally investigates the analytic results reported in [5], [25], [10] and [33]. [11] considers sizing buffers managed with active queues management techniques.

The foregoing work is in the context of wired links, and to our knowledge the question of buffer sizing for 802.11 wireless links has received almost no attention in the literature. Exceptions include [21] [23] [29]. Sizing of buffers for voice traffic in WLANs is investigated in [21]. The impact of fixed buffer sizes on TCP flows is studied in [23]. In [29], TCP performance with a variety of AP buffer sizes and 802.11e parameter settings is investigated. In [16] [17], initial investigations are reported related to the eBDP algorithm and the ALT algorithm of the A* algorithm. We substantially extend the previous work in this paper with theoretical analysis, experiment implementations in both testbed and a production WLAN, and additional NS simulations.

VIII. CONCLUSIONS

We consider the sizing of network buffers in 802.11 based wireless networks. Wireless networks face a number of fundamental issues that do not arise in wired networks. We demonstrate that the use of fixed size buffers in 802.11 networks inevitably leads to either undesirable channel underutilization or unnecessary high delays. We present two novel buffer sizing algorithms that achieve high throughput while maintaining low delay across a wide range of network conditions. Experimental measurements demonstrate the utility of the proposed algorithms in a real environment with real traffic.

The source code used in the NS-2 simulations and the experimental implementation in MadWifi can be downloaded from www.hamilton.ie/tianji_li/buffersizing.html.

APPENDIX

In the production WLAN of the Hamilton Institute, the AP is equipped with an Atheros 802.11/a/b/g PCI card and an external antenna. The operating system is a recent Fedora 8 (kernel version 2.6.24.5). The latest MadWifi driver version (0.9.4) is used, in which the buffer size is fixed at 250 packets. The AP is running in 802.11g mode with the default rate adaptation algorithm enabled (i.e., SampleRate [7]). All data traffic is processed via the Best Effort queue, i.e., MadWifi is operating in 802.11 rather than 802.11e mode. A mix of Windows/Apple MAC/Linux laptops and PCs use the WLAN from time to time.

REFERENCES

[1] IEEE 802.11 WG. International standard for information technology – local and metropolitan area networks, part 11: wireless LAN MAC and PHY specifications, 1999.

[2] *Part 11: wireless LAN medium access control (MAC) and physical layer (PHY) specifications: Medium Access Control (MAC) Quality of Service (QoS) Enhancements*, IEEE 802.11e/D8.0, Feb. 2004.

[3] S. A. Mujtaba, et. al., "TGn Sync Proposal Technical Specification," www.tgnsync.org, IEEE 802.11-04/889r6, May 2005.

[4] Online, madwifi-project.org.

[5] G. Appenzeller, I. Keslassy, and N. McKeown, "Sizing Router Buffers," in *Proc. of ACM SIGCOMM*, 2004, pp. 281–292.

[6] N. Beheshti, Y. Ganjali, M. Ghobadi, N. McKeown, and G. Salmon, "Experimental Study of Router Buffer Sizing," in *Proc. of IMC*, Oct. 2008.

[7] J. Bicket, "Bit-rate Selection in Wireless Networks," MSc. Thesis, MIT 2005.

[8] C. Chatfield, *The Analysis of Time Series, An Introduction*, CRC Press 2004.

[9] A. Dhamdher and C. Dovrolis, "Open Issues in Router Buffer Sizing," in *Computer Communication Review*, Jan. 2006.

[10] M. Enachescu, Y. Ganjali, A. Goel, N. McKeown, and T. Roughgarden, "Routers with Very Small Buffers," in *Proc. of INFOCOM*, Dec. 2006.

[11] D. Y. Eun and X. Wang, "Achieving 100% Throughput in TCP/AQM Under Aggressive Packet Marking With Small Buffer," *IEEE/ACM Transactions on Networking*, vol. 16, no. 4, pp. 945 - 956, Aug. 2008.

[12] C. Kellett, R. Shorten, and D. Leith, "Sizing Internet Router Buffers, Active Queue Management, and the Lur'e Problem," in *Proc. of IEEE CDC*, 2006.

[13] K. Kumaran, M. Mandjes, A.L. Stolyar, "Convexity Properties of Loss and Overflow Functions," *Operations Research Letters*, vol. 31, no.2, pp. 95-100, 2003.

[14] A. Lakshminantha, R. Srikant, and C. Beck, "Impact of File Arrivals and Departures on Buffer Sizing in Core Routers," in *Proc. of Infocom*, Apr. 2008.

[15] D. Leith, P. Clifford, D. Malone, and A. Ng, "TCP Fairness in 802.11e WLANs," *IEEE Communications Letters*, vol. 9, no. 11, Jun. 2005.

[16] T. Li and D. Leith, "Buffer Sizing for TCP Flows in 802.11e WLANs," *IEEE Communications Letters*, Mar. 2008.

[17] T. Li and D. Leith, "Adaptive Buffer Sizing for TCP Flows in 802.11e WLANs," *Chinacom 2008*.

[18] T. Li, Q. Ni, D. Malone, D. Leith, T. Tuletto, and Y. Xiao, "Aggregation with Fragment Retransmission for Very High-Speed WLANs," *IEEE/ACM Transactions on Networking*, vol. 17, no. 2, pp. 591-604, Apr. 2009.

[19] D. Malone, K. Duffy, and D.J. Leith, "Modeling the 802.11 distributed coordination function in non-saturated heterogeneous conditions," *IEEE/ACM Transactions on Networking*, vol. 15, no. 1, Feb. 2007.

[20] D. Malone, D.J. Leith, A. Aggarwal, and I. Dangerfield, "Spurious TCP Timeouts in 802.11 Networks," *WiNMe*, Apr. 2008.

[21] D. Malone, P. Clifford, and D. J. Leith, "On Buffer Sizing for Voice in 802.11 WLANs," *IEEE Communications Letters*, vol. 10, no. 10, pp 701–703, Oct. 2006.

[22] V. Paxson and S. Floyd, "Wide-Area Traffic: The Failure of Poisson Modeling," *IEEE/ACM Trans. Netw.*, vol. 3, no. 3, pp. 226-244, Jun. 1995.

[23] S. Pilosof, et. al., "Understanding TCP fairness over Wireless LAN," in *Proc. of IEEE INFOCOM 2003*.

[24] R. Prasad, C. Dovrolis, and M. Thottan, "Router Buffer Sizing Revisited: The role of the input/output capacity ratio," *IEEE/ACM Transactions on Networking*, to appear.

[25] G. Raina and D. Wischik, "Buffer Sizes for Large Multiplexers: TCP Queueing Theory and Instability Analysis," in *Proc. of EuroNGI*, Jul. 2005.

[26] R. Shorten, F. Wirth, and D. Leith, "A Positive Systems Model of TCP-Like Congestion Control: Asymptotic Results," *IEEE/ACM Transactions on Networking*, vol. 14, no. 3, pp. 616–629, Jun. 2006.

[27] R. Stanojevic, C. Kellett, and R. Shorten, "Adaptive Tuning of Drop-Tail Buffers for Reducing Queueing Delays," *IEEE Communications Letters*, vol. 10, no. 7, pp 570–572, Jul. 2006.

[28] D. Tang and M. Baker, "Analysis of A Local-Area Wireless Network," in *Proc. of ACM MobiCom*, Aug. 2000.

[29] M. Thottan, and M. C. Weigle, "Impact of 802.11e EDCA on mixed TCP-based applications," in *Proc. of IEEE WICON 2006*.

[30] O. Tickoo and B. Sikdar, "On the Impact of IEEE 802.11 MAC on Traffic Characteristics," *IEEE J. on Selected Areas in Commun.*, vol. 21, no. 2, Feb. 2003, pp. 189-203.

[31] C. Villamizar and C. Song, "High Performance TCP in ANSNET," *ACM Computer Communication Review*, vol. 24, no. 5, pp. 45–60, Oct. 1994.

[32] G. Vu-Brugier, R. Stanojevic, D. Leith, and R. Shorten, "A Critique of Recently Proposed Buffer-Sizing Strategies," *ACM Computer Communication Review*, vol. 37, no. 1, Jan. 2007.

[33] D. Wischik and N. McKeown, "Part I: buffer sizes for core router," *ACM Computer Communication Review*, vol. 35, no. 3, Jul. 2005.

[34] Y. Zhang and D. Loguinov, "ABS: Adaptive Buffer Sizing for Heterogeneous Networks," in *Proc. of IWQoS*, Jun. 2008.

[35] Z. Zhao, S. Darbha, and A. L. N. Reddy, "A Method for Estimating the Proportion of Nonresponsive Traffic At a Router," *IEEE/ACM Transactions on Networking*, vol. 12, no. 4, pp. 708–718, Aug. 2004.

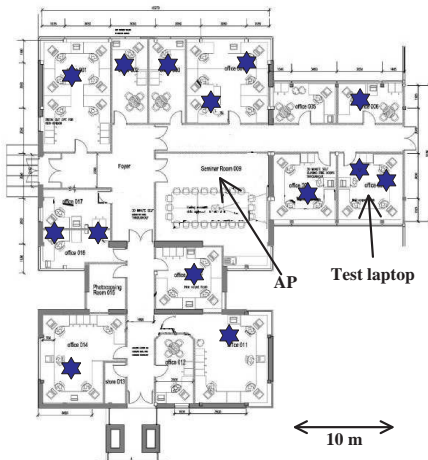


Fig. 23. WLAN of the Hamilton Institute. Stars represent users' approximate locations.

PLACE PHOTO HERE

Tianji Li received the M.Sc. (2004) degree in networking and distributed computation from École Doctorale STIC, Université de Nice-Sophia Antipolis, France, and the Ph.D. (2008) degree from the Hamilton Institute, National University of Ireland Maynooth, Ireland, where he is currently a research fellow. He is interested in improving performance for computer and telecommunication networks.

PLACE PHOTO HERE

Douglas Leith graduated from the University of Glasgow in 1986 and was awarded his PhD, also from the University of Glasgow, in 1989. In 2001, Prof. Leith moved to the National University of Ireland, Maynooth to assume the position of SFI Principal Investigator and to establish the Hamilton Institute (www.hamilton.ie) of which he is Director. His current research interests include the analysis and design of network congestion control and distributed resource allocation in wireless networks.

PLACE PHOTO HERE

David Malone received B.A.(mod), M.Sc. and Ph.D. degrees in mathematics from Trinity College Dublin. During his time as a postgraduate, he became a member of the FreeBSD development team. He is a research fellow at Hamilton Institute, NUI Maynooth, working on wireless networking. His interests include wavelets, mathematics of networks, IPv6 and systems administration. He is a co-author of O'Reilly's "IPv6 Network Administration".