

Active Queue Management for Fair Bandwidth Allocation of Mixed Responsive and Non-Responsive Traffic Using a Closed-Loop Congestion Control Scheme

Roman Pletka[†], Andreas Kind[†], Marcel Waldvogel[†] and Soenke Mannel[‡]

[†]IBM Zurich Research Laboratory, CH-8803 Rüschlikon, Switzerland

{rap, ank, mwl}@zurich.ibm.com

[‡]ISR, Universität Stuttgart, D-70511 Stuttgart, Germany

soenke.mannel@isr.uni-stuttgart.de

Abstract—Today’s known and widely used active queue management (AQM) schemes do not differentiate between packets from responsive (e.g., TCP sessions) and non-responsive traffic (e.g., UDP). This results in further widening the gap of unfair advantage already inherent to non-responsive traffic, as the responsive sender will significantly reduce its future transmit rate as a result of the congestion signals. As a simple work-around, responsive and non-responsive traffic are often assigned distinct AQM parameters. This approach however requires tuning for each traffic class that potentially depends on the current or expected offered load. In other words, responsiveness and TCP-friendliness cannot be estimated easily—not at last due to short-lived TCP sessions. In this paper we propose a closed-loop congestion control (CLCC) scheme on top of an existing AQM scheme to achieve fair bandwidth distribution among concurrent responsive and non-responsive traffic. The new scheme has the advantage that it does not need to estimate the level of responsiveness of traffic. We analyze our scheme on top of an existing rate-based AQM scheme known to approximate max-min fairness, and by means of simulations show that our extension significantly improves fair bandwidth allocation for responsive and non-responsive traffic.

I. INTRODUCTION

AQM systems control packet-drop probabilities according to the level of congestion. Typically, packets are dropped with a higher probability when congestion increases in order to avoid undesired lock-out situations and unnecessarily high delays [1]. Past research efforts focused on configuring AQM systems such as RED [2] and numerous variants thereof [3]–[7] for service-level specifications. In particular, the difficulty in finding the appropriate parameters for RED for any combination of offered load was pointed out in [8], [9].

In general, AQM schemes are not able to achieve accurate bandwidth guarantees. Depending on the underlying hardware, a full-featured scheduler might not be available (e.g., there is only a priority scheduler, a WRR scheduler, or a FIFO queue), and therefore absolute bandwidth guarantees are difficult to achieve solely through AQM. Even in the presence of a separate fair queuing (FQ) scheduler, the goals of AQM and FQ will partially disagree. Namely, AQM tries to keep queues short, while FQ needs longer, fixed-length queues to perform its task. Under congestion, FQ will fill its buffers, preventing them from accepting bursts and adding to the total queuing delay.

As a result, AQM is often beneficiary to FQ. Especially rate-based AQM [10] seems to be more convenient for Quality-of-Service (QoS) and Service Level Agreement (SLA) environments, as it allows configuration based on packet rates rather than buffer thresholds. However, as will be shown later, even

with rate-based AQM, some manual tuning is necessary to achieve fair bandwidth allocation for responsive (e.g., TCP) and non-responsive (e.g., UDP) traffic at the same time. This is due to end-to-end congestion control mechanisms built into responsive protocols whereby packet drops are interpreted as network congestion somewhere on the transit path causing the sending rate to be reduced (back-off mechanism). This congestion response makes responsive protocols willfully give up bandwidth to avoid the dreaded congestion collapse. Non-responsive, greedy protocols on the other hand will grab the new bandwidth. Given sufficient bandwidth requested by non-responsive traffic, these will continuously reduce the bandwidth available to responsive, “nice” traffic, until the latter are left without any bandwidth.

In this paper we propose a new closed-loop congestion control (CLCC) scheme on top of an AQM scheme that enables fair bandwidth allocation even under heavily bursty traffic conditions with concurrent responsive and non-responsive networking traffic. In particular, we use Bandwidth Allocation Technology (BAT) [10] as the underlying AQM algorithm to show the benefit of our control scheme. We expect the new CLCC scheme to work just as easily with other types of AQM schemes, especially with those based on additive increase, multiplicative decrease (AIMD) control algorithms.

AIMD schemes require the configuration of at least two main parameters per traffic class, namely the increase and decrease constants. Thus in a setup with two traffic classes, i.e., a responsive and a non-responsive one, four parameters need to be tuned. Our new control scheme adapts and optimizes only one of these four BAT parameter that before used to be constant. Extensive simulations showed that the new CLCC scheme significantly outperforms the traditional BAT algorithm, and is effective when executed on a much slower time scale than drop probability updates with BAT are. Therefore the scheme needs not necessarily be implemented directly in hardware. However, an implementation in hardware should be considered if the amount of control traffic is an issue.

The paper is structured as follows. First, we briefly discuss related work, motivating the need for an improved AQM scheme. Then, in Section III we first introduce and evaluate the BAT algorithm in presence of responsive and non-responsive traffic. Section IV proposes our new scheme, CLCC, that improves fairness characteristics in this situation. Results obtained from various simulations are presented in Section V. We conclude in Section VI.

II. RELATED WORK

Dynamic-RED [11] uses control theory for adapting the packet-dropping probability to stabilize the queue length close to a given threshold value. Benefits are bounded delays and independence from the number of flows present in the system. As the algorithm is not rate-based, it cannot provide bandwidth guarantees.

BLUE [12] has been introduced to overcome the problems encountered by AQM schemes based on queue lengths. The algorithm uses packet-loss and link-idle events to manage congestion, and the proposed extension called *Stochastic Fair BLUE* is able to identify and rate-limit non-responsive flows by using multiple levels of independent hash functions. The probability that flows are misclassified and therefore needlessly discriminated increases with the number of flows present in the system. A major drawback is the additional expense in implementation complexity.

The properties of congestion feedback mechanisms based on AIMD and their fairness properties were discussed extensively [13]–[16]. BAT by Bowen *et al.* [10] applies this background towards a rate-controlled AQM scheme. This proves effective at keeping the queue length to a minimum while absorbing short-term bursts. BAT approximates max-min fairness [17] by managing the drop rates using an open-loop control. Individual offered loads act as a feedback signal. Compared to other AQM algorithms, BAT also performs well for non-responsive networking traffic.

III. OPEN-LOOP CONTROL FOR BANDWIDTH ALLOCATION

A. The BAT Algorithm

The BAT active queue-management scheme [10], an open-loop control system, maintains per-flow transmit probabilities $T_i(t)$ and operates on two (coarse and fine grained) AIMD levels. For a flow i , the probability function is given as:

$$T_i(t + \Delta t) = \begin{cases} \min(1, T_i(t) + w) & \text{if } f_i(t) \leq f_{i,\min} \\ T_i(t)(1 - w) & \text{else if } f_i(t) > f_{i,\max} \\ \min(1, T_i(t) + C_i \cdot \bar{B}(t)) & \text{else if } B(t) = 1 \\ T_i(t)(1 - D_i \cdot O_i(t)) & \text{otherwise} \end{cases} \quad (1)$$

The more aggressive part of the algorithm brings the actual flow rate into the maximum and minimum bounds ($f_{i,\min}, f_{i,\max}$) configured, while a fine-tuned control loop operates within these bounds. The aggressive part is governed by the constant w ($0 < w < 1$). The coefficients C_i and D_i determine the rate of increase and decrease during fine-grained control. Their choice depends on the desired convergence rate. The two-level technique of BAT increases its convergence speed while retaining the stability properties of AIMD [13].

At initialization, the transmit probability for each flow i is $T_i(0) = 1$, and T_i are then computed periodically in intervals of Δt from the preceding transmit probability value, the current offered rate of the flow O_i , the current serviced rate of the flow given as $f_i = T_i \cdot O_i$, and the exponentially weighted moving average $\bar{B}(t)$ of the excess bandwidth signal $B(t)$. $B(t)$ indicates whether transmit bandwidth is available or is expected to become available soon [10].

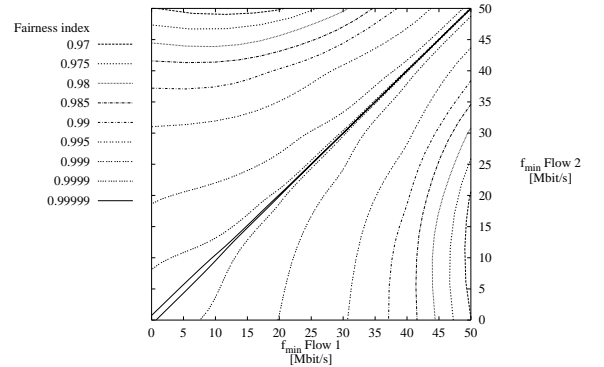


Fig. 1. Fairness index for two flows depending on their $f_{i,\min}$ value. $f_{i,\max}$ is set to 100 Mbit/s for both flows.

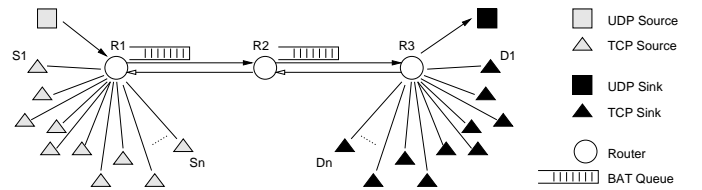


Fig. 2. Simulation topology for measuring combined UDP and TCP traffic using BAT.

B. Fairness of Non-Responsive Traffic

First we investigate the fairness of AQM in presence of non-responsive traffic only. To do so, we have chosen the BAT algorithm as our reference AQM. To illustrate the fairness of BAT, we analyze two competing CBR sources sending UDP traffic at 100 Mbit/s over a 100-Mbit/s link while varying their minimum bandwidth guarantee $f_{i,\min}$ from 0 to 50 Mbit/s. In our implementation, BAT updates transmit fractions every 4 ms. The average bandwidth allocated to each flow over a simulation time of $t = 20$ s is compared with the theoretical fair value of max-min fairness. As a metric for fairness, we use a weighted version of the fairness index as introduced in [18]:

$$f(x_1, x_2, \dots, x_n) = \left(\sum_{i=1}^n \frac{x_i}{w_i} \right)^2 / n \sum_{i=1}^n \left(\frac{x_i}{w_i} \right)^2. \quad (2)$$

The fairness index lies between 0 and 1, and is 1 if all users receive the theoretical throughput, given by the weights w_i . Thus, Equation 2 tells how close to a desired bandwidth allocation, in our case max-min fairness, the measurements are. The simulation results in Figure 1 show the weighted fairness index as a function of various minimum bandwidth guarantees. High fairness is achieved as long as the $f_{i,\min}$ are close and decreases whenever they drift away from each other as well as for higher minimum bandwidth guarantees.

C. Responsive vs. Non-Responsive Traffic

The results of [10] show that BAT performs well for non-responsive traffic and keeps the buffer occupancy level low during periods of heavy congestion. However, fair bandwidth allocation is typically not achieved in presence of *both* responsive *and* non-responsive traffic because the congestion control used by responsive protocols amplifies the control signal.

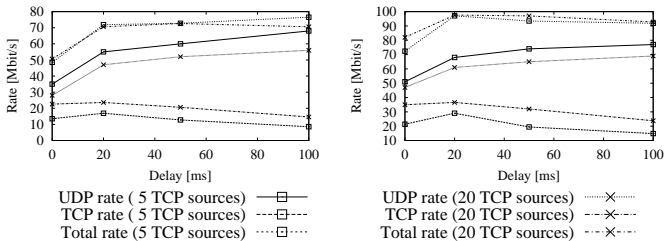


Fig. 3. Measured TCP throughput, allocated UDP bandwidth, and total throughput for different f_{\min} values of BAT. f_{\max} is set to 100 Mbit/s.

Figure 2 shows the network topology used for our simulations, which closely matches the hardware testbed with BAT implemented on an IBM PowerNP 4GS3 network processor. Traffic is divided into two flows: The first consists of a non-responsive CBR flow, and the second of responsive TCP connections. The CBR flow is generated using an IXIA traffic generator; TCP sources reside on a Linux computer with a 2.4.18 kernel using `ttcp` [19]. On the receiving side, the network emulator DummyNet [20] delays the packets before they are delivered to the application (`ttcp` receiver) to simulate WAN distances.

Figure 3 shows the UDP and TCP rates allocated by BAT for different $f_{i,\min}$ as a function of the delay introduced for various numbers of TCP sources. The allocated TCP bandwidth increases with the number of TCP sources. This is consistent with findings from other AQM algorithms [11], [21] and can be modeled using the throughput equation of TCP [22]. Furthermore, the TCP bandwidth share decreases with increasing RTT, as it will take longer to recover from packet loss [22].

D. A First Step to Improve Fairness

A first measure to increase fairness between responsive and non-responsive traffic would be to adapt the increase and decrease constants C_i and D_i depending on the traffic characteristics currently observed. Thus, for a responsive and a non-responsive traffic, four parameters need to be adjusted. In early simulations we found that it is sufficient to adapt the decrease constant of non-responsive traffic D_{UDP} to achieve significantly better fairness results. Higher values for D_{UDP} allow a faster reaction to non-responsive traffic in the case of congestion, compensating for the non-responsive characteristics of UDP traffic.

Using *ns* [23] and the simulation topology in Figure 2, we first investigate the fairness behavior of BAT. The link rate for all links is set to 100 Mbit/s. A high number of TCP sources fed by greedy FTP agents share a bottleneck link with non-responsive UDP traffic created by a CBR source. All TCP connections have different RTTs that are equally distributed within the interval given in Table I. Clearly, because of the congestion control mechanism of TCP, the two flows will not get the same bandwidth allocated. With identical increase and decrease constants C_i and D_i , non-responsive traffic is able to grab an almost 10 times higher bandwidth share.

Figure 4(a) shows the linearly interpolated relation between the allocated bandwidth and D_{UDP} gained from a large set

TABLE I

Simulation	UDP rate [Mbit/s]	RTT [ms]	Packet size [bytes]
0–4	100	2–20	1500
5–9	200	2–20	1500
10–14	100	2–20	500
15–19	200	2–20	500
20–24	100	2–198	1500
25–29	200	2–198	1500
30–34	100	2–198	500
35–39	200	2–198	500

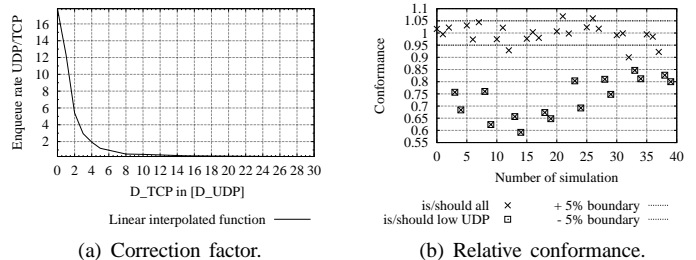


Fig. 4. Correction factor for D_{UDP} as a function of the desired bandwidth allocation given by maximum and minimum bandwidth guarantees obtained by a linear interpolation over a large set of different simulations described in Table I, and its relative conformance with the set of simulations.

of simulations. Detailed simulation parameters for the open-loop case are listed in Table I. Each simulation set consists of five different bandwidth allocations, namely 90, 60, 45, 30 and 25% UDP traffic, respectively. For the set of simulations given, conformance often exceeds 5% (Figure 4(b)), especially when the allocated UDP rate is relatively small (indicated with crossed boxes) compared with the offered TCP load. In summary, the method introduced above fails to provide reasonable long-term fairness.

IV. IMPROVING FAIRNESS WITH CLCC

To improve fairness, we introduce a closed-loop control mechanism. The goal of the control process is to allocate bandwidth in a fair way as a long-term objective. When choosing the actual controller, we started out with the family of PID (Proportional-integral-derivative) controllers [24], which are characterized by relatively low computational requirements, suitability for hardware implementation, and their excellent theoretical performance.

An ideal PID controller is described as follows:

$$y(t) = K_R \left[x_d(t) + \frac{1}{T_N} \int x_d(t) dt + T_V \frac{dx_d(t)}{dt} \right], \quad (3)$$

where the time-dependent variables are the control input $y(t)$, and the control error $x_d(t)$. K_R , T_N and T_V are system parameters. The terms in the square brackets represent the proportional, integral, and differential influence (from left to right). The discrete time notation of Equation 3 can be found by approximation of the derivatives using backward differences. Solving for the actual value of the reference input yields

$$y_k = y_{k-1} + b_0 \cdot x_{d,k} + b_1 \cdot x_{d,k-1} + b_2 \cdot x_{d,k-2}, \quad (4)$$

where b_0 , b_1 , and b_2 determine the basic properties of the controller. By setting $b_1 = b_2 = 0$, Equation 4 corresponds

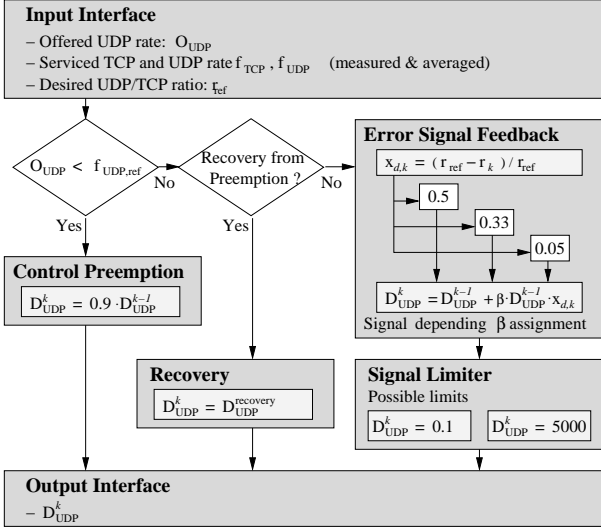


Fig. 5. Diagram of the AQM control algorithm using CLCC.

to a pure integral controller. This simplification leads to an additional reduction in both complexity and calculation effort while still ensuring efficient elimination of constant steady-state errors. We found that a proportional part ($b_1 \neq 0$) does not lead to a clear improvement of the results. Therefore we limit ourselves to integral feedback control.

The main idea of our CLCC is to adapt only the non-responsive decrease factor D_{UDP} , which in the traditional BAT model was a predetermined constant [10]. The now time-dependent value of D_{UDP} is determined by a control process that adjusts its value autonomously according to the current traffic characteristics. This decision was made because responsive protocols react strongly to congestion indications; therefore, we need to adapt the behavior of non-responsive traffic forcibly by tuning D_{UDP} .

The controlled variable is the ratio of serviced UDP and TCP rates r_k . The reference variable is represented by r_{ref} and corresponds to the desired ratio derived from max-min fairness [17]. The control error $x_{d,k}$ is defined as

$$x_{d,k} = (r_{\text{ref}} - r_k) / r_{\text{ref}}. \quad (5)$$

Now D_{UDP} is a time-dependent value specifying the manipulated variable of the system given as

$$D_{\text{UDP}}^k = D_{\text{UDP}}^{k-1} + b_0 \cdot \frac{r_{\text{ref}} - r_k - 1}{r_{\text{ref}}}, \quad (6)$$

which describes the difference equation of a time-discrete, strictly bilinear system of first order.

To adapt the control-error signal to the current state of D_{UDP} , we define $b_0 = \beta \cdot D_{\text{UDP}}^{k-1}$; $\beta < 0$. The parameter β depends on the absolute value of $x_{d,k}$. Three different cases for β are distinguished that stress the importance of $x_{d,k}$ for large deviations, as detailed in Figure 5. The ratio between UDP and TCP traffic r_k is measured every 100 ms and averaged over a time interval of 1 s to obtain a low-pass filter effect. The control algorithm is invoked every second while the underlying AQM scheme calculates new transmit fractions every 4 ms.

The detailed CLCC algorithm is shown in Figure 5. The ratio r_k is measured every 100 ms and averaged over a time interval of 1 s to obtain a low-pass filter effect. The control algorithm is invoked every second, whereas the underlying AQM scheme calculates new transmit fractions every 4 ms.

To increase stability, we bound D_{UDP} as shown in Figure 5. Moreover, the control loop is suspended when the offered UDP load falls below $f_{\text{UDP,ref}}$, because a useful adaptation is no longer possible. The system is then said to be in control-preemption state, in which the preceding value of D_{UDP} is preserved.

Next, we present results from simulations done in *ns* and verified with a prototype implementation on a IBM PowerNP 4GS3 network processor.

V. RESULTS

In Section III-D, we showed that D_{UDP} should not be a constant value. As a consequence thereof, we introduced CLCC. Here we compare simulations of the BAT AQM schemes with and without CLCC. We use identical simulation setups and take the default increase and decrease parameters for the open-loop case. The exact simulation setup is as follows: One UDP flow is sending at a constant bitrate of 90 Mbit/s and shares a 100-Mbit/s bottleneck link between router R1 and R2 with a set of 100 greedy TCP connections, each of them having a uniformly distributed RTT between 2 and 200 ms. The controller allocates link bandwidth in the ratio of $r_{\text{ref}} = 30/70 \approx 0.42$. The results are shown in Figure 6.

In Figure 6(a), the unacceptable difference between desired and actual ratio using traditional BAT for responsive traffic is evident. This behavior is also seen in the spread of the fairness index. The improvements obtained using our algorithm are shown in Figure 6(b). Although the UDP/TCP ratio still fluctuates, the long-term average closely matches the desired value, which is also reflected in improvements in the fairness index. Further simulations with Pareto-distributed on/off-sources for TCP traffic showed similar results, and are not presented here due to space limitations.

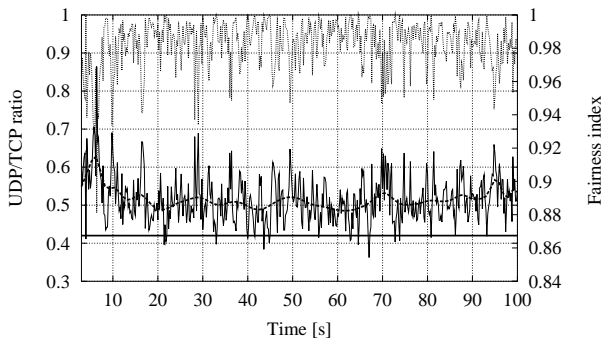
In a next step, we tested CLCC in presence of a UDP load alternating between 50 and 120 MBit/s every five seconds. r_{ref} is set to 0.4. The results of this simulation are presented in Figure 7.

Clearly, the new control algorithm is also suitable in presence of bursty traffic. Fair bandwidth allocation is rapidly recovered even after large changes in the offered UDP rate.

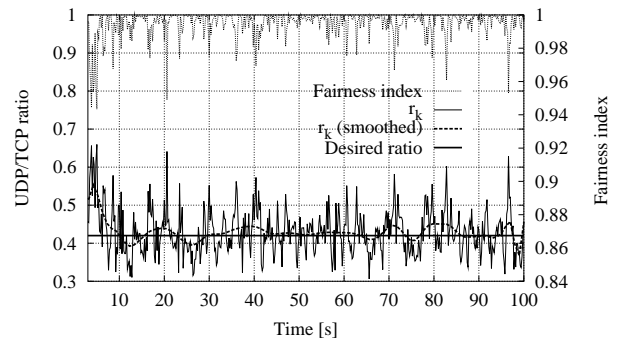
VI. CONCLUSION

In this paper, the problem of fair bandwidth allocation for mixed responsive and non-responsive traffic is addressed. We proposed a closed-loop congestion control scheme for an existing AQM system and showed that the fairness characteristics can be improved significantly. The control scheme is based on a proportional integral derivative controller and has been simulated in *ns* and implemented on an IBM PowerNP 4GS3 network processor for verification of the results presented.

We argue that AQM with intrinsic fairness properties such as BAT with CLCC provides advantages over combinations of other AQM schemes combined with an FQ scheduler. Specifically, the separate buffering required by FQ will partly

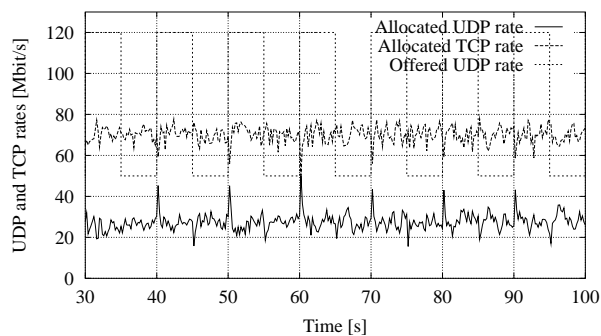


(a) BAT without additional controller.

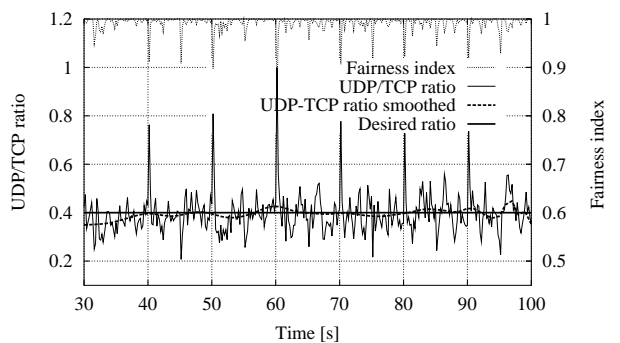


(b) Using the CLCC algorithm.

Fig. 6. Simulation results with a set of greedy TCP connections and one constant UDP flow.



(a) UDP and TCP rates.



(b) UDP/TCP ratio r_k and fairness index.

Fig. 7. Simulation results with a set of 100 greedy TCP connections and one bouncing UDP flow.

negate the short-queue goal of AQM. To absorb bursts, FQ queue lengths need to be set to constant, high values. Under heavy load, this leads to full buffers introducing delays and preventing the queues from accepting bursts.

To our knowledge, this is the first AQM system that is able to automatically control under disparate reactivity without having to also include an expensive fair queuing system.

REFERENCES

- [1] B. Braden *et al.* Recommendations on queue management and congestion avoidance in the Internet. RFC 2309, Internet Engineering Task Force, April 1998.
- [2] S. Floyd and V. Jacobson. Random early detection gateways for congestion avoidance. *ACM Trans. on Networking*, August 1993.
- [3] D. Lin and R. Morris. Dynamics of random early detection. In *Proc. SIGCOMM '97*, pages 127–137, Cannes, France, September 1997.
- [4] D. Clark and W. Fang. Explicit allocation of best effort packet delivery service. *ACM Trans. on Networking*, 6(4), August 1998.
- [5] V. Rosolen, O. Bonaventure, and G. Leduc. A RED discard strategy for ATM networks and its performance evaluation with TCP/IP traffic. *ACM Computer Commun. Rev.*, 29(3), July 1999.
- [6] W. Feng, D. D. Kandlur, D. Saha, and K. G. Shin. A self-configuring RED gateway. In *Proc. of INFOCOM '99*, pages 1320–1328, March 1999.
- [7] Cisco Systems Online IOS Documentation. Weighted Random Early Detection (WRED), 1998.
- [8] M. May, J. Bolot, C. Diot, and B. Lyles. Reasons not to deploy RED. In *Proc. of 7th Int. Workshop on Quality of Service (IWQoS '99)*, London, pages 260–262, June 1999.
- [9] V. Misra, W.-B. Gong, and D. F. Towsley. Fluid-based analysis of a network of AQM routers supporting TCP flows with an application to RED. In *Proc. SIGCOMM 2000*, pages 151–160, August 2000.
- [10] E. Bowen, C. Jeffries, L. Kencl, A. Kind, and R. Pletka. Bandwidth allocation for non-responsive flows with active queue management. In *Proc. Int. Zurich Seminar on Broadband Commun., IZS 2002*, February 2002.
- [11] J. Aweya, M. Ouellette, D. Y. Montuno, and A. Chapman. A control theoretic approach to active queue management. *Computer Networks*, 36:203–235, 2001.
- [12] W. Feng, D. Kandlur, D. Saha, and Kang G. Shin. BLUE: A new class of active queue management algorithms. Technical Report CSE-TR-387-99, University of Michigan, April 1999.
- [13] D. Chiu and R. Jain. Analysis of the increase and decrease algorithms for congestion avoidance in computer networks. *Computer Networks & ISDN Systems*, 17(1):1–14, June 1989.
- [14] F. P. Kelly and A. Maulloo. Rate control in communication networks: Shadow prices, proportional fairness and stability. *Operational Res. Soc.*, (49):237–252, 1998.
- [15] P. Hurley, J.-Y. Le Boudec, and P. Thiran. A note on the fairness of additive increase and multiplicative decrease. In *Proc. ITC-16, Edinburgh, UK*, pages 467–478, June 1999.
- [16] S. Gorinsky and H. Vin. Additive increase appears inferior. Technical Report 2000-18, Dept. of Computer Sciences, Univ. of Texas, May 2000.
- [17] D. Bertsekas and R. Gallager. *Data Networks*. Prentice-Hall, 1987.
- [18] R. Jain. *The Art of Computer Systems Performance Analysis*. John Wiley, New York, 1991.
- [19] M. Muuss. The story of the test TCP program ttcp. <http://ftp.arl.army.mil/~mike/ttcp.html>.
- [20] Luigi Rizzo. Dummynet: A simple approach to the evaluation of network protocols. *ACM Computer Commun. Rev.*, 27(1):31–41, 1997.
- [21] R. Morris. TCP behavior with many flows. In *Proceedings of IEEE International Conference on Network Protocols (ICNP)*, pages 205–211, October 1997.
- [22] J. Padhye, V. Firoiu, D. Towsley, and J. Kurose. Modeling TCP throughput: A simple model and its empirical validation. In *Proc. SIGCOMM*, pages 303–314, 1998.
- [23] S. Bajaj *et al.* Improving simulation for network research. Technical Report 99-702b, Univ. Southern California, March 1999.
- [24] G. F. Franklin, J. D. Powell, and M. Workman. *Digital Control of Dynamic Systems*. Addison-Wesley, third edition, 1997.