

# Phantom: A Simple and Effective Flow Control Scheme

Yehuda Afek      Yishay Mansour      Zvi Ostfeld

Computer Science Department,  
Tel-Aviv University, Israel 69978.  
{afek,mansour,ostfeld}@math.tau.ac.il

## Abstract

This paper presents *Phantom*, a simple constant space algorithm for rate based flow control. As shown by our simulations, it converges fast to a fair rate allocation while generating a moderate queue length. While our approach can be easily implemented in ATM switches for managing ABR traffic, it is also suitable for flow control in TCP router based networks. Both the introduced overhead and the required modifications in TCP flow control systems are minimal. The implementation of this approach in TCP guarantees fairness and provides a unifying interconnection between TCP routers and ATM networks. The new algorithm easily inter-operates with current TCP flow control mechanisms and thus can be gradually introduced into installed based TCP networks.

## 1 Introduction

A major goal of network protocols is to maximize the utilization of the network resources (e.g., bandwidth) while sharing the resources in a fair way among the users/applications. Flow control is the mechanism that controls the flow of traffic into the network and whose goal is to avoid and resolve data traffic congestion while ensuring high utilization and fairness among the different connections. As the speed and size of computer networks increase the flow control is becoming more and more critical and challenging. In such networks a small mistake by the flow control for a tiny period of time may quickly result in clogging and in the loss of many messages.

Communication protocols in general and in large high speed networks in particular should have the following properties: simplicity, minimal space requirements (both buffers and algorithm variables), short response time to bursty traffic requirements, and inter-operability across several networks. The flow control mechanism is no exception; in the flow control algorithm presented here has all the above essential features.

One property, fairness, is unique to flow control. While

intuitively we may understand the notion of fairness, namely, fairly sharing the network bandwidth among different users, defining it rigorously in a network environment is not trivial. A widely acceptable measure of fairness is the *max-min fairness* [BG87, Jaf81, GB84, Gaf82, Hay81, Mos84, Cha94]. A bandwidth allocation is max-min fair if for every session, one cannot increase its bandwidth without decreasing the bandwidth of sessions of equal or lower bandwidth [BG87].

This paper presents an efficient, fair and simple scheme for flow control, called *Phantom*, which is suitable for both TCP router based networks and ATM networks. Thus, not only suggesting a flow control mechanism for the two network types, but also making the first step in integrating the flow control of both types of networks.

The key idea behind the Phantom algorithm is to bound the rate of sessions that share a link by the amount of unused bandwidth on that link. As though there is another imaginary session, a *phantom* session, on the link, and the link bandwidth is fairly shared among all the sessions, including the phantom. Using the Phantom mechanism max-min fairness is achieved for a network in which on every link there is an extra “phantom” session.

In addition the Phantom scheme has the following advantages: (1) it requires a constant number of variables in every output port of any router/switch, (2) it quickly and robustly responds to variations in the traffic, (3) the scheme utilizes any amount of available bandwidth, and can be easily integrated in a network where the amount of available bandwidth may change over time, since the scheme does not use explicitly the specific bandwidth allocation, (4) in certain variations of Phantom the forward RM cells are not necessary. (5) the scheme may be implemented with simple hardware, using only basic arithmetics, without multiplications or divisions.

The idea of adding another imaginary session to every link is not completely new, Jaffe [Jaf81] used a similar idea in order to study the convergence of flow control algorithms. However, the algorithm of Jaffe [Jaf81], as well as other related works along this lines [GB84, Hay81, Mos84], may not be effectively used in a distributed environment in which the traffic frequently changes.

Several flow control mechanisms have been suggested in recent years for TCP router based networks [Jac88, BP95, Zha89, FJ93]. In [Jac88, BP95] the flow control operates at the source ends of the sessions by adjusting the window size according to the detection of packet loss

or the round trip delay. Other mechanisms that operate at the routers (gateways) are Random Early Detection [Zha89, FJ93] and Early Packet Discard [RF94], that in case of congestion drops new incoming packets in an intelligent approach.

Different ways of applying the Phantom scheme to flow control in TCP routers based networks are presented in Section 3. While the main advantage of these applications is in achieving fairness, it does require some modification of either the source end, or the routers, and the TCP packet format. However, it seems that fairness may not be achieved without some modifications either at the routers [Kes91] or at the TCP header.

Applying the Phantom scheme to ATM networks is somewhat easier due to the handles provided by the ATM Forum. ABR (Available Bit Rate) is a class of service in ATM networks that allows bursty data applications to utilize the available network bandwidth in a *fair* way, while guaranteeing a low cell loss rate. The ATM Forum on Traffic Management has adopted rate based as the basic mechanism for flow control of ABR traffic. The basic principle of the rate based flow control is as follows<sup>1</sup>: A control message (Resource Management cell, RM cell) is looping around the virtual circuit of each session. On its forward way the RM cell carries the CCR (Current Cell Rate) of the session and “informs” the switches about the session’s rate ACR (allowed Cell Rate), while on the backward direction, each switch may decrease the ER (Explicit Rate) field in the RM cell, according to the congestion that it observes. When the RM cell arrives back to the source, the value in the ER field is the maximum rate at which the source may transmit. However, the source may not increase its rate too fast, the amount of increase is limited by a parameter, AIR (Additive Increase Rate) per RM cell received [Sat96].

Since the ATM Forum has adopted the rate based scheme as a standard for flow control, [BF95], several proposals for its implementation in a switch have been proposed. One way of partitioning the different proposals is by their space requirements: (1) Flow control algorithms with constant space, i.e., the number of variables used by the algorithm is constant independent of the number of sessions passing through a port [Rob94, ST94, JKV94, Bar94, Bar95], and (2) unbounded space algorithms, i.e., algorithms in which the space depends (usually linearly) on the number of connections [CCJ95, KVR95, CR96, TW96, JKG<sup>+</sup>95]. The Phantom algorithm presented here is a constant space algorithm. In Section 4 it is compared to four other constant space algorithms that have been presented in the ATM Forum. We believe it overcomes some of the weaknesses of these algorithms.

Constant space flow control algorithms are important for today’s networks, e.g., if one would like to support both ATM and TCP traffic, then a dependence on the number of sessions is prohibitory.

## 2 Phantom - The Basic Scheme

The basic idea of the Phantom algorithm is to keep a certain portion of the link capacity unused and to limit the rates of sessions sharing the link by the amount of

<sup>1</sup>Following the DECbit [RJ90], the ATM flow control supports another mechanism using the CI (Congestion Indication) bit on control cells

the *unused* bandwidth on that link (as if there is an extra phantom session on the link). Specifically, we define  $\Delta$  to be the unused link capacity, i.e.,  $Link\ Capacity - \sum rates\ of\ sessions\ that\ use\ the\ link$ . The rate of sessions that are above  $\Delta$  is reduced towards  $\Delta$  and the rate of sessions that are below may be increased. This mechanism reaches a steady state only when the unused capacity ( $\Delta$ ) is equal to the maximum rate of any session that crosses the link and all the sessions that are constrained by this link are at this rate. The value of  $\Delta$  is easily computed in each output port by counting the number of cells arriving at the queue of that port over an interval of time and subtracting this amount from the number of cells that may be transmitted in that interval (this quantity is normalized by dividing it by the length of the time interval).

For example, if three sessions share a 100 Mbs link, then in steady state each session receives 25 Mbs and the link utilization is at 75% (in this case  $\Delta$  is 25 Mbs too). However, if two of the three sessions are restricted elsewhere to 10 Mbs each, the third session gets 40 Mbs and  $\Delta$  is also 40 Mbs.

The basic scheme of Phantom is very robust, as our simulations indicate. The robustness stems from the following observations:

1. In a stable environment the scheme converges to a steady state in which fairness is achieved. This fact follows from the property that all the sessions that are constrained on some link get the same rate (which is the value of  $\Delta$  for that link) in a stable state.
2. Since the Phantom scheme preserves a residual unused capacity, it smoothly accommodates the addition of a new session without a queue buildup.
3. The basic attitude of Phantom is to **avoid** congestion rather than to **react** upon it. For example, consider a 100 Mbs link shared by three sessions, two of which are constrained elsewhere to 10 Mbs, and one limited by this link to 40 Mbs. If the 10 Mbs restriction of the two sessions is removed at the same time then while the rate of the two sessions increases the rate of the third session decreases, even before a queue is built. This is because the value of  $\Delta$  decreases as the rate of the two sessions increases above 10 Mbs.<sup>2</sup>

However, to employ the Phantom scheme in an actual network tuning and adjustments are necessary. Since the necessary mechanisms for its implementation in ATM are readily available we start by describing and analyzing it in an ATM environment. However, we believe the scheme is attractive and effective also in TCP networks and that the ATM analysis and results would basically carryover to the TCP settings. In Section 3 we discuss its implementation for flow control over TCP.

### 2.1 Simulations configuration in ATM

Our simulations are done in BONEs [ALT94]. The ATM end systems are as specified in [Sat96], Appendix I. The

<sup>2</sup>Assuming that the rate of the two sessions does not instantly jump to a very high value, that is, in ATM terms [Sat96],  $AIR \cdot Nrm$  is much smaller than 30 Mbs.

major parameters of the end systems are as follows (unless specified otherwise):  $Nrm = 32$ ,  $AIR \cdot Nrm = 42.5Mbs$ ,  $RDF = 256$ ,  $PCR = 150Mbs$ ,  $TOF = 2$  and  $ICR = 8.5Mbs$ . Though we use the optional use-it-or-lose-it behavior, the use of this option by the sources does not effect the correctness of the algorithm (and in particular, it does not effect most of our simulations, where the sources are greedy and the "lose-it" behavior is unlikely to show-up). The links parameters (bandwidth and delay) are given separately for each simulation. The parameters used in Phantom are specified in the sequel.

## 2.2 Implementation in ATM

The following five issues, concerning the implementation of the algorithm, are addressed in the sequel:

1. **Measuring  $\Delta$ :** A naive measurement of  $\Delta$  may be very noisy and unstable due to the inherent delay in the system and to the nature of the traffic. Hence, a technique to filter out the variations and noise in the measurement is necessary. The filtered value of  $\Delta$  is called *MACR* (details are given in the sequel).
2. **Sensitivity to queue length:** Inaccuracies, rounding errors, and noise, may still introduce an error which might accumulate over time and exhibit itself in a gradual queue buildup. Hence, the operation of the algorithm has to be coupled with the absolute value of the queue length.
3. **Utilization:** If only a few "heavy" sessions use the network then the scheme might under-utilize the available bandwidth. Below we introduce a slight modification to ensure high utilization also in these cases.
4. **Variance Consideration:** In some circumstances a minor change in the value of *MACR* may cause a large surge in the incoming traffic. This in turn may result in large oscillations in *MACR* even if the sources are greedy (steady). We introduce a method that reduces these oscillations, while still maintaining the fast response to changes in the sources' traffic.
5. **Reducing the maximum queue length:** Under certain scenarios, e.g., many synchronized sessions starting to transmit at the same time, a large queue may momentarily develop. Here we present a method that allows us to tune the trade off between the size of such a spike in the queue length and the response time of the algorithm.

**Measuring  $\Delta$ :** The amount of residual bandwidth,  $\Delta$ , is measured in fixed time intervals of length  $\tau$  and is accumulated in a parameter called Maximum Allowed Cell Rate (MACR) [Rob94, Bar95]. *MACR*, is updated at the end of each time interval by the weighted sum of its previous value and the measured  $\Delta$ . I.e.,:

$$MACR = MACR \cdot (1 - \alpha) + \Delta \cdot \alpha, \quad \alpha \in [0, 1]$$

In our simulations  $\tau$  is set to 100 cell time and  $\alpha$  is set to  $1/16$ .

While this method seems to work reasonably it turns out that a more robust mechanism is necessary. Consider

the following scenario: a sudden burst causes  $\Delta$  to be significantly smaller (or even have a negative value, if the traffic destined to that output port exceeds the link capacity) for a short period of time. This sharp change in  $\Delta$  causes *MACR* to be significantly smaller which in turn caused the sessions to reduce their rates, and hence the network to be under utilized. Now, because a low rate session transmits fewer RM cells, it takes the system a long time to gain back the normal utilization. To this end, and to avoid the possibility of erroneous negative *MACR*, we prevent *MACR* from dropping sharply, i.e.,:

$$MACR := \max(MACR \cdot (1 - \alpha) + \Delta \cdot \alpha, MACR \cdot dec\_factor), \\ dec\_factor \in [0, 1]$$

(In our simulations  $dec\_factor = 3/4$ .)

In Fig. 1 the Phantom algorithm as described so far, is employed to control the rate of four sessions sharing one 150 Mbs link. The figure shows *MACR*, the queue length at the output port as well as the sessions' ACR.

Fig. 2 presents the behavior of the algorithm under the same scenario but the four sessions have a diverse range of RTTs.

In Fig. 3 the fast reaction of the scheme to the addition of a fifth session is shown.

In Fig. 4 the algorithm is examined with on/off traffic.

**Sensitivity to queue length:** Since the algorithm is insensitive to the absolute queue length it might happen that the queue length is growing without bound. Such a problematic case is likely to happen when many sessions are restricted by the same link. This phenomena is not unique to this algorithm, but is shared by many algorithms that are insensitive to the queue length. Fig. 5 illustrates such a scenario.

In order to solve this problem we extend the algorithm by a simple addition that makes it sensitive to the queue length. Instead of using a fixed value  $\alpha$  in the computation of *MACR* we use two different variables  $\alpha_{inc}$  and  $\alpha_{dec}$ . The variable  $\alpha_{inc}$  is used instead of  $\alpha$  when  $\Delta > MACR$  and  $\alpha_{dec}$  is used when  $\Delta \leq MACR$ . Moreover, the actual value of  $\alpha$  depends on the queue length, when the queue length is relatively small  $\alpha_{inc}$  is relatively large and  $\alpha_{dec}$  is small. On the other hand, when the queue length is large,  $\alpha_{dec}$  is relatively large and  $\alpha_{inc}$  is relatively small. This method shortens the convergence time of sessions and thus decreases the period of time in which the link is either underutilized or overutilized.

The values of  $\alpha_{inc}$  and  $\alpha_{dec}$  used in our simulations (Fig. 5) are presented in Fig. 6 as a function of the current queue length and the parameters *QT* (Queue Threshold, which is 50 cells in our simulations).

**Utilization:** The underutilization of the algorithm when only very few sessions share a link, may be avoided by using the following modification to the basic scheme. Rather than restricting the sessions by  $\Delta$  we may restrict them by a constant (called *utilization\_factor*) times *MACR*. For example, if we restrict the sessions by  $3 \cdot MACR$ , and we have three sessions sharing a 100Mbs link, then each session is allocated 30 Mbs (and the phantom session gets 10 Mbs). Our simulations show that this modification does not effect the other characteristics of our scheme. (The only possible differences are larger transient queue length

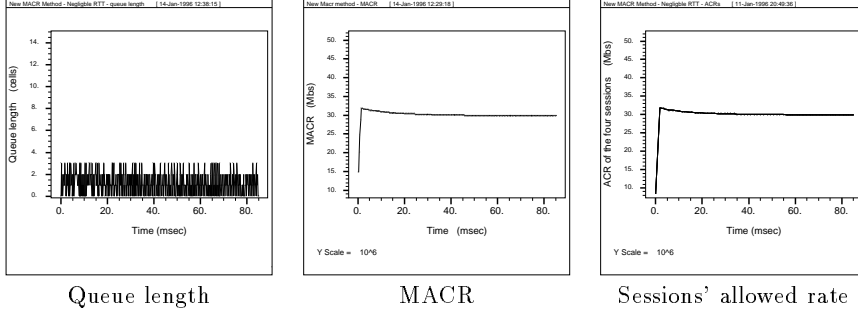


Figure 1: Phantom - Four greedy sessions with negligible RTT (0.01 ms) sharing a 150 Mbs link. Both  $MACR$  and the sessions' rate converge to  $\frac{150}{4+1} = 30Mbs$ .

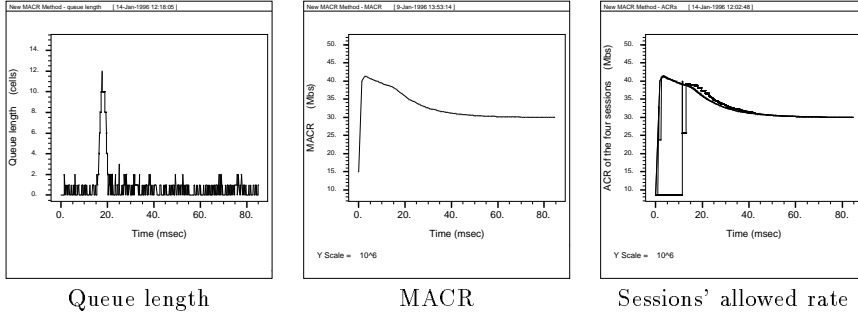


Figure 2: Phantom - Four greedy sessions as in Fig. 1, with different RTTs (11 ms, 1 ms, 0.1 ms and 0.01 ms) sharing a 150 Mbs link. Three sessions with a relatively small RTT get high rates at the beginning (they quickly receive the feedback from the switch) while the fourth session still uses its initial rate. However, it does not take long before the fourth session gets the feedback and converges to the desired value (which is  $\frac{150}{4+1} = 30Mbs$ ). Note that a small queue (12 cells) is formed during convergence time and is drained shortly after.

queue length		$\alpha_{inc}$	$\alpha_{dec}$
queue length	$< QT$	$\alpha$	$\alpha/16$
$QT <$	queue length $\leq 2QT$	$\frac{\alpha}{2}$	$\alpha/8$
$2QT <$	queue length $\leq 4QT$	$\frac{\alpha}{4}$	$\alpha/4$
$4QT <$	queue length $\leq 8QT$	$\frac{\alpha}{8}$	$\alpha/2$
$8QT <$	queue length $\leq 16QT$	$\frac{\alpha}{16}$	$\alpha$
$16QT <$	queue length	$\frac{\alpha}{32}$	$2 \cdot \alpha$

Figure 6: Phantom - Computing  $\alpha_{inc}$  and  $\alpha_{dec}$  according to the current queue length.

and larger oscillations in the value of  $MACR$ .) In Fig. 7 the simulation of Fig. 1 is repeated but the sessions rates are restricted by  $5 \cdot MACR$  instead of just  $MACR$ .

**Variance Consideration:** When many sessions share one link with a high utilization factor there may be large oscillations in  $MACR$  and in the queue length (Fig. 8). To eliminate this phenomena we first approximate the standard deviation in  $\Delta$ , and then take it into consideration in the calculation of  $\alpha_{inc}$  and  $\alpha_{dec}$ .

Following [Jac88] the standard deviation is approximated by the mean deviation (which is easier, as it does not require a square root computation) as follows:

$$ERR := MACR - \Delta$$

$$\sigma := \sigma \cdot (1 - h) + |ERR| \cdot h, \quad h \in [0, 1]$$

(In our simulations  $h$  is set to  $1/16$ ).

However, using this computation one can not distinguish between the case in which  $\sigma$  is large due to an external change, such as the addition or removal of a new session, and the case in which the large variation is caused by a small change in the calculation of  $MACR$  that is then amplified by the large gain in the feedback loop (many sessions with large utilization factor). Note that each of the two cases requires a different treatment. In the first case Phantom should normally react with the large gain when updating  $MACR$  (adapting to an external change in the traffic) while in the second case, the gain in the control loop should be reduced to flatten the oscillations in  $MACR$ .

To distinguish between the two cases we use two additional variables  $\sigma_{pos}$  and  $\sigma_{neg}$ . When  $\Delta \geq MACR$   $\sigma_{pos}$  is used and when  $\Delta \leq MACR$   $\sigma_{neg}$  is used. The variable  $\sigma$  is then set to be the minimum between  $\sigma_{pos}$  and  $\sigma_{neg}$ . The pseudo code for the estimation of  $\sigma$  is as follows:

**Every  $\tau$  seconds time interval do:**  
 $ERR := MACR - \Delta$   
 If  $ERR \geq 0$  then  
 $\sigma_{neg} := \sigma_{neg} \cdot (1 - h) + ERR \cdot h$   
 Else  
 $\sigma_{pos} := \sigma_{pos} \cdot (1 - h) + |ERR| \cdot h$   
 $\sigma := \min(\sigma_{neg}, \sigma_{pos})$   
 $ratio := f(\sigma)$  { See tabular of  $f$  below. }  
 $\alpha_{inc}, \alpha_{dec} := g(queue\ length, \alpha)$  { See Fig. 6 }  
 $\alpha_{inc} := \alpha_{inc} \cdot ratio$  ;  $\alpha_{dec} := \alpha_{dec} \cdot ratio$

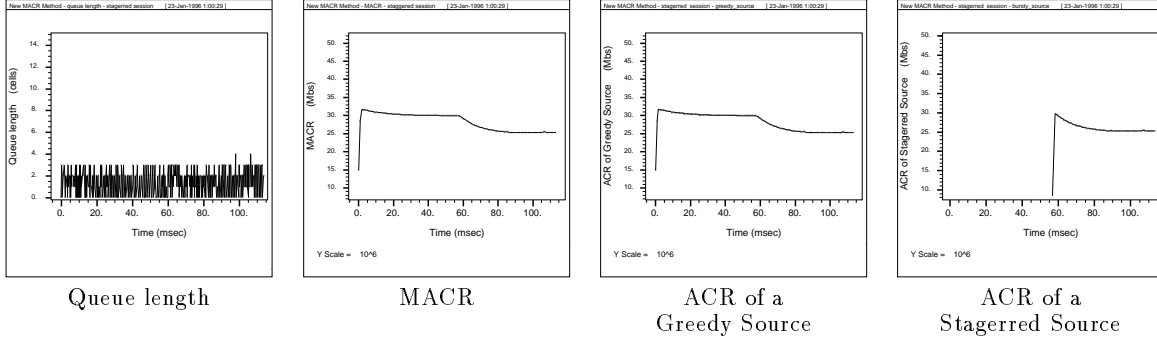


Figure 3: Phantom - Four greedy sessions sharing a 150 Mbs link; The allocation for each session is  $\frac{150}{4+1} = 30 Mbs$ . A staggered session is added after 58 ms. The addition does not cause a queue buildup and all the sessions' rate converge to a new value of  $\frac{150}{5+1} = 25 Mbs$ .

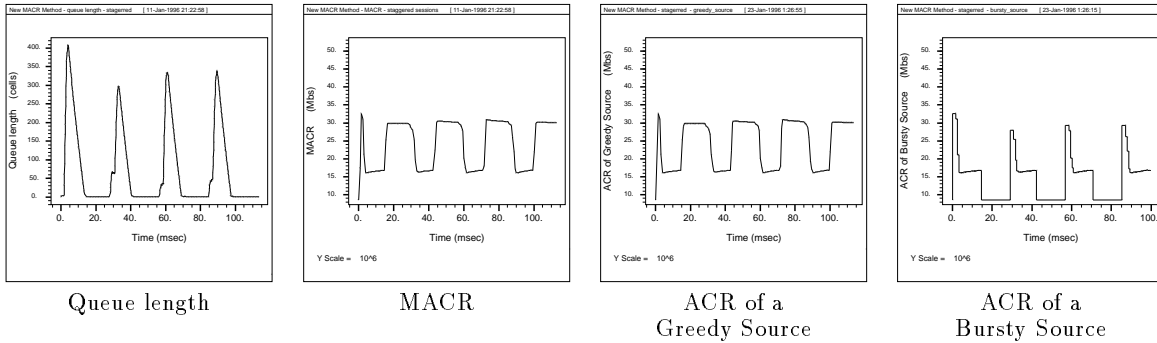


Figure 4: Phantom - Four greedy sessions and four on/off sessions sharing a 150 Mbs link. The four on/off sessions are on for a period of 14 ms and then become idle for a period of 14 ms.  $MACR$  is either  $\frac{150}{4+1} = 30 Mbs$  when there are 4 active sessions or  $\frac{150}{8+1} = 16.7 Mbs$  when there are 8 active sessions. Note that a queue is built when the on/off sessions become active and is drained out afterwards. When the on/off sessions become idle adjustment is achieved without queue buildup. Note also that the  $ACR$  of the on/off sessions is set to be their Initial Cell Rate ( $ICR$ ), 8.5 Mbs when they are idle. However, during that period no data is transmitted by those sessions.

The motivation for using such an approach is the following: In the case where the large variance is caused by an external change only one of the variables  $\sigma_{pos}, \sigma_{neg}$  gets a large value. However, if the large variance is caused by perpetual big oscillations then both variables are relatively big. The value of  $ratio$  may be computed from the value of  $\sigma$  as follows:

$\sigma$	$ratio := f(\sigma)$
$\sigma \leq MACR$	1
$MACR < \sigma \leq 2 MACR$	1/2
$2 MACR < \sigma \leq 4 MACR$	1/4
$4 MACR < \sigma \leq 8 MACR$	1/8
$8 MACR < \sigma$	1/16

The right hand diagrams of Fig. 8 present the behavior of the algorithm with the variance method. We have also tested the behavior of the algorithm in a dynamic changing environment with four greedy sessions and four on/off sessions in an environment identical to that of Fig. 4 and the results were exactly the same as in Fig. 4. Thus testifying that the variance method does not seem to affect the behavior of the algorithm when there are large external changes in the traffic.

**Reducing the maximum queue length:** Consider the following scenario: The  $MACR$  of a link that is shared by many sessions is suddenly increased from a small value to a relatively large value (because e.g., some VBR connection stopped). In this case all the sessions that were bottlenecked on this link increase their rates to  $MACR \cdot utilization\_factor$  at about the same time (assume large AIR), resulting in a surge of traffic much larger than the link capacity thus developing a large queue (see Fig. 9).

We employ two mechanisms to significantly reduce this large developing queue. First, (following the ideas of the TCP window flow control algorithm,) rather than allowing a session rate to jump at once to  $MACR \cdot utilization\_factor$  we restrict its rate increase to at most twice its current rate during each round trip time (by allowing the ER on the backward RM cell to be at most twice the value of the CCR on the RM cell). The price of this mechanism is a somewhat slower convergence when only a few sessions share the link. Fig. 10 examines this new attitude in the same environment as the simulation of Fig. 9.

While the first method has significantly reduced the maximum queue length a careful observation of Fig. 10 reveals that the rate of an arbitrary session is still increasing when there is already a significant queue buildup

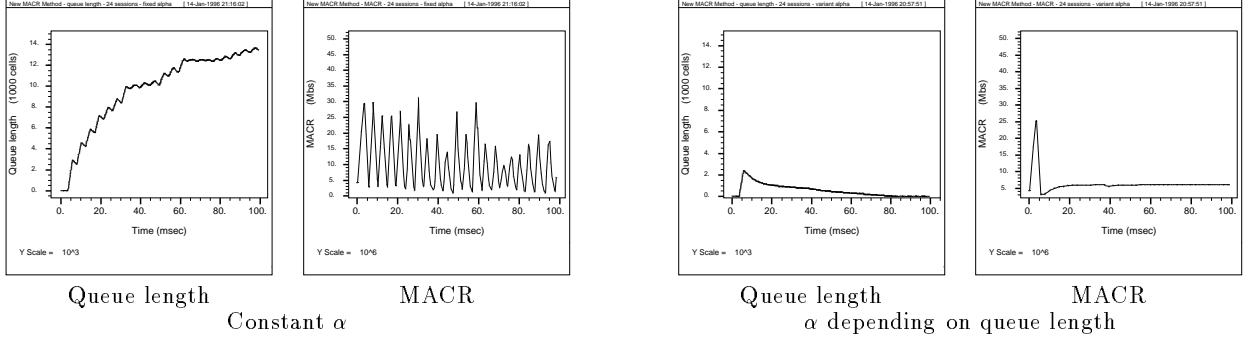


Figure 5: Phantom - Twenty four greedy sources sharing a 150 Mbs link, each with a negligible RTT (0.01 ms). The two left hand diagrams illustrate the problematic case where  $\alpha$  is constant. There is a "synchronization" effect between the sessions, i.e., whenever there is a small change in the value of  $MACR$ , it causes a large total change in the rate of all the sessions together. This change causes a large change in  $\Delta$  and vice versa. Hence, there are large oscillations of the sessions rates which do not converge even after a long period of time and the queue length grows without bound. In the two right hand diagrams  $\alpha_{inc}$  and  $\alpha_{dec}$  are computed according to Fig. 6. The queue length peaks at 2,000 rather than continuously increasing, and  $MACR$  converges nicely and smoothly to the desired value (which is  $\frac{150}{24+1} = 6Mbs$ ). The large burst at the beginning of the simulation is caused by the "synchronized" burst of 24 sources. However, after this initial peek the value of  $MACR$  is stable and the queue is drained.

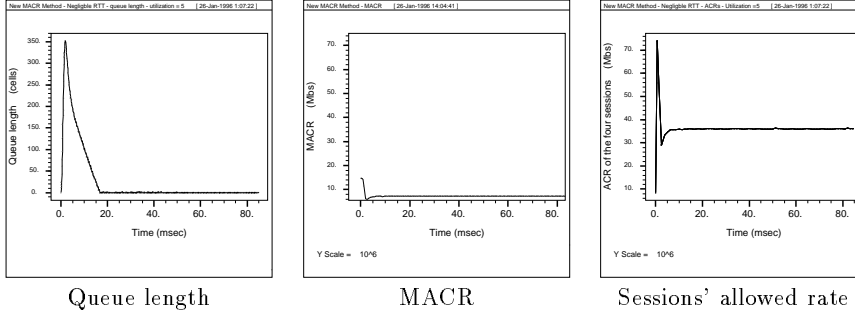


Figure 7: Phantom - Four greedy sessions with negligible RTT sharing a 150 Mbs link. The ratio between  $MACR$  and the link restriction is 5. Each session stabilizes on a rate of  $\frac{150}{4+\frac{1}{5}} = 35.7Mbs$  while  $\Delta$  stabilizes on a value of 7.14 Mbs. Hence, the utilization during steady state in this case is more than 95%. However, one can see that during convergence time, the queue length might have a temporary larger value (350 cells), compared with that in Fig. 1 (where no queue was built up).

(over utilization). Therefore we employ a second mechanism that monitors the available bandwidth more closely than  $MACR$ . The new measure is called  $Fast\_MACR$  and is computed in the same way as  $MACR$ , i.e.:

$$Fast\_MACR := Fast\_MACR \cdot (1 - \beta) + \beta \cdot \Delta, \quad \beta \in [0, 1]$$

where  $\beta$  is larger than  $\alpha$  ( $\beta = 1/8$  in our simulations), thus  $Fast\_MACR$  is more sensitive to fast changes in the available link bandwidth. Whenever  $MACR$  is significantly larger than  $Fast\_MACR$  we assume that  $MACR$  is too large and set the  $NI$  (No Increase) bit in each backward  $RM$  cell. Any source that observes this bit set, may not increase its rate (on top of all the other restrictions, [Sat96]). Fig. 11 illustrates the effect of this method on the same scenario as in Fig. 9.

Fig. 12 is a detailed version of the pseudo code of the algorithm.

### 2.3 Convergence Analysis

$FS$ , the *Fair Share* of a link is the maximum possible value such that the rate of every session that crosses the link is

either restricted to be exactly  $FS$  or restricted elsewhere to a lower value. Let  $ls$  be the set of sessions whose rate is smaller than  $FS$ ,  $\sum ls$  the total rate of sessions in  $ls$ ,  $\#s$  the total number of sessions sharing the link and  $\#ls$  the number of sessions in  $ls$  and  $C$  be the link bandwidth, then  $FS = \frac{C - \sum ls}{\#s - \#ls}$ . Note that a set of flows is max-min fair if the flow of each session equals the minimum fair share along its path [BG87, Jaf81, GB84, Cha94, Jai95].

Using the above terms we define for the Phantom algorithm,  $FS_{Phantom}$  for every link to be

$$\frac{C - \sum ls}{\#s - \#ls + \frac{1}{k}}$$

where  $k$  is the *utilization\_factor* used in the algorithm.

Clearly when the rate of every session is the minimum  $FS_{Phantom}$  along its path, the network is in a stable state (assuming steady traffic requirements). Note that if  $k = 1$  then the fair share of *Phantom* is exactly the same as that of the max-min allocation, when on every link there is an additional imaginary greedy session. As  $k$  is increased the bandwidth consumed by the phantom session in a stable state is decreased though the allocation between the other

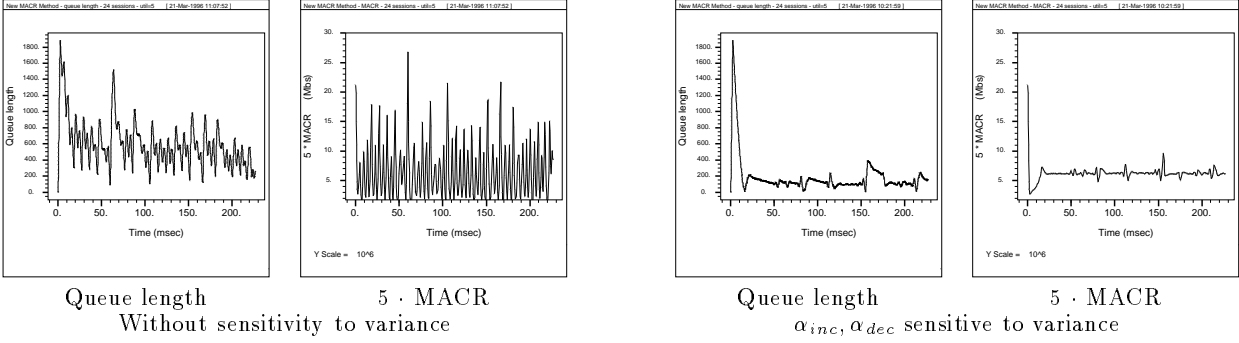


Figure 8: Phantom - twenty four greedy sessions with negligible RTT sharing a 150 Mbs link where *utilization\_factor* = 5. The reason for the large oscillations in  $\Delta$  is the fact that a small change in the value of *MACR* is multiplied by 5 and affects 24 sessions (i.e., the optimum value of *MACR* is 1/121 of the link capacity). The two right hand diagrams demonstrate how the variance method overcomes this problem by drastically reducing the oscillations in *MACR* and thus also reducing the queue length.

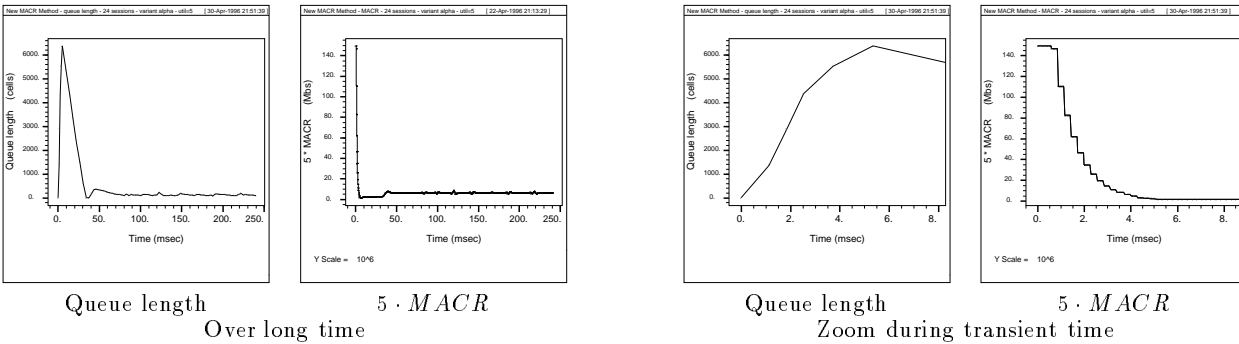


Figure 9: Phantom - 24 greedy sessions with negligible RTT sharing a 150 Mbs link and *utilization\_factor* = 5. The initial value of  $5 \cdot MACR$  is the full link capacity, i.e., 150 Mbs. The initial rate of each session (*ICR*) is relatively small (1.7 Mbs). The large queue length (more than 6000 cells) is created during the transient time, until *MACR* becomes close to its optimal value (6.2 Mbs).

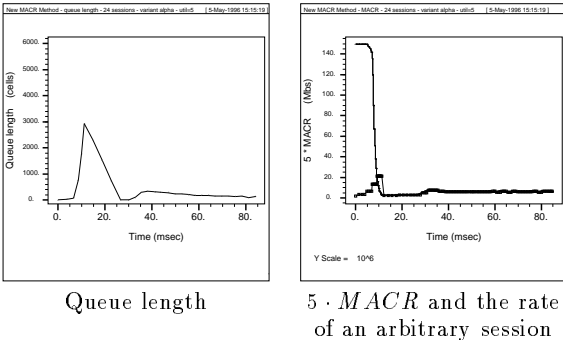


Figure 10: Phantom - 24 greedy sessions with negligible RTT sharing a 150 Mbs link where *utilization\_factor* = 5. Sessions may at most double their rate each round trip time. The maximum queue length is now 3000 cells.

sessions is still *fair*. The convergence of our algorithm to the steady state follows from arguments similar to those in [Cha94, KVR95].

We have tested our algorithm in other configurations such as the Generic Fairness Configuration suggested in [Woj94]. Due to space limitations these results are not

included but can be found in [AMO96].

#### 2.4 Additional Remarks:

1. Phantom is set to acquire any available bandwidth as soon as it becomes available. Because, in calculating  $\Delta$  we may subtract the total incoming traffic (from all traffic classes) from the total link capacity thus acquiring immediately any bandwidth that becomes available (the UBR traffic may be excluded depending on the relative priority).
2. In stable states the Phantom value of *MACR* is the maximum rate permitted for a session. This value may be instrumental for traffic policing. E.g., if the value of *CCR* of a session is correct and represents its actual rate then, misbehavior of a session is easily identified (*CCR* verification is easier).

### 3 Applying Phantom to flow control in TCP

#### 3.1 TCP flow control techniques

In the past few years several new techniques for flow control of TCP traffic in IP networks were introduced and

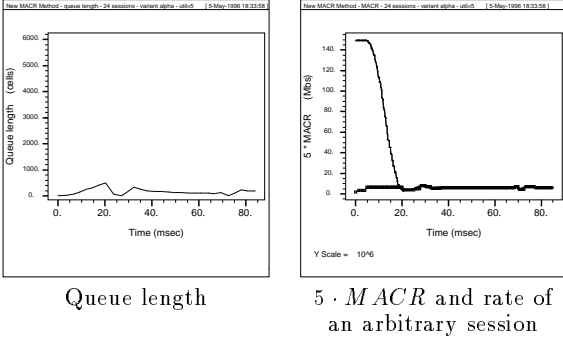


Figure 11: Phantom - 24 greedy sessions with negligible RTT sharing a 150 Mbs link where  $utilization\_factor = 5$ . Sessions' rate is at most doubled each round trip time and the *NI* bit is set if *MACR* is significantly larger than *Fast\_MACR*. The maximum queue length is now 500 cells. Note that as soon as congestion starts, the increase in the rate of the selected session stops.

employed. We distinguish between flow control mechanisms that operate at the end stations and those that operate at the routers/gateways.

**Source end stations mechanisms:** Two such algorithms are currently proposed, Reno, suggested by Jacobson [Jac88], and Vegas, by Brakmo and Paterson [BP95]. Both algorithms operate by dynamically adjusting the TCP window size. If the network is observed not congested the window size is increased while if the network is observed congested, the window size is decreased. The basic idea is to consider the network congested if packet loss is detected. In Vegas there are several additional mechanisms one of which indicates congestion if the round trip delay (the time interval between the departure of a packet and the arrival of its acknowledgment) exceeds some threshold.

**Router mechanisms:** Obviously when the router buffer overflows some packets are dropped. Another option for the router is to send an ICMP Source Quench message [Ste94, BP87], which causes the source to decrease its window size. This option however, is rarely used as the Source Quench message may overload an already congested network. Furthermore, this mechanism is biased against certain kinds of traffic as shown in [FJ92].

A more sophisticated strategy, called RED (Random Early Detection) was presented by Floyd and Jacobson [FJ93], where either packets are dropped or Source Quench messages are sent when the router's queue length exceeds a certain threshold. The packet to be dropped (or the source to which a Source Quench message is sent) is selected according to some probability distribution on the packets. The RED approach reduces the queue lengths at the router and overcomes some of the bias described in [FJ92]. Yet, the resulting mechanism still does not always guarantee fairness (as discussed in the sequel).

For networks where each router sends packets in a round robin order between the sessions, another router mechanism was presented by Keshav [Kes91]. In this method the source measures the maximum queuing delay by sending two consecutive packets and measuring the time difference at which they return to the source. This

**Constant:**  
 $C$  { The total link capacity }

**Parameters:**  
 $\tau, \alpha, \beta, QT, decreasing\_factor, h, utilization\_factor$

**Variables:**  
 $MACR, \Delta, \alpha_{inc}, \alpha_{dec}, \sum data$

**Initialization:**  
 Initialize the values of *MACR* and *Fast\_MACR* to be some portion of  $C$ .  
 $\sum data := 0$

**When a forward data is received:**  
 $\sum data := \sum data + size\ of\ current\ data$

**Every  $\tau$  seconds time interval do:**  
 Choose  $\alpha_{inc}$  and  $\alpha_{dec}$  according to the current queue length and the current deviation.  

$$\Delta := \min(C - \frac{\sum data}{\tau}, \frac{C}{utilization\_factor})$$
 If  $\Delta > MACR$  then  

$$MACR := MACR \cdot (1 - \alpha_{inc}) + \Delta \cdot \alpha_{inc}$$
 Else  

$$MACR := \max(MACR \cdot (1 - \alpha_{dec}) + \Delta \cdot \alpha_{dec}, MACR \cdot decreasing\_factor)$$

$$Fast\_MACR := Fast\_MACR \cdot (1 - \beta) + \Delta \cdot \beta$$
 $\sum data := 0$

**For every backward RM cell do:**  
 $ER\ field\ on\ cell := \min(utilization\_factor \cdot MACR, current\ ER\ field\ on\ cell, 2 \cdot CCR\ field\ on\ cell)$   
 If  $(MACR > 2 \cdot Fast\_MACR)$  then  
 Set the value of the *NI* bit in the *RM* cell

Figure 12: Phantom - Detailed Pseudo Code.

mechanism may achieve high degree of fairness and does not depend on any specific network topology, however its implementation complicates the router architecture (requires a queue per session) and requires the cooperation of all the routers.

For further discussion on routers congestion control see the survey of Mankin and Ramakrishnan [MR91].

Three problems are likely to occur in TCP traffic flow control that does not use explicit rate indication:

**Unfair allocation:** With the exception of [Kes91] the above mechanisms do not guarantee fairness. For example when two sources that use Vegas [BP95], got different window sizes, and both have the same delay thresholds ( $\alpha, \beta$ ) then, there is no mechanism that would balance them. The current mechanisms would either increase both or decrease both. Another source for unfairness in allocation is the bias against sessions that pass through many routers (analogous to the "beat down" phenomena in ATM, [BdJ94]). The two diagrams on the left hand side of Fig. 13 and Fig. 16 demonstrate an unfair behavior of Reno in an environment of drop tail routers.

**Congestion detection:** Some of the above mechanisms base their estimate of the queuing delay on the knowledge of the round trip delay (RTT). However, RTT is the sum of the queuing delays and propagation delays along the path. As the source uses only the total delay it is almost

impossible to separate the two components, leaving a large degree of uncertainty about the queuing delay.

**Sensitivity to parameters:** Some of the above schemes are sensitive to variations in parameters values. For example, out of tune time threshold parameters may cause network underutilization (see the discussion in [BP95] for example). Alternatively, distinct parameters for different sessions may cause severe unfairness. E.g., two sessions using Vegas and sharing one router such that the lower time threshold ( $\alpha$ ) of the one is larger than the upper time threshold ( $\beta$ ) of the other. Then convergence is achieved only when the first session uses the entire link bandwidth.

### 3.2 Employing Phantom in TCP

Here we show how the employment of the Phantom approach may solve some of the problems described above. We suggest and discuss different possibilities of applying the basic Phantom scheme to the flow control of TCP traffic over IP networks. Our suggestions introduce minimal overhead, while achieving smooth interaction with current TCP equipment, thus it can be gradually introduced into TCP environments. For the implementation of Phantom, TCP routers would have to compute  $MACR$ , either as explained above or by measuring the amount of unutilized link bandwidth (which is an approximation for the value of  $MACR$ ).

There are two basic approaches in the application of Phantom to TCP:

**Without modifying the routers:** The source periodically polls the routers on the path to the destination using e.g., SNMP. In each polling the source fetches the  $MACR$  values of those routers, and adjusts its window according to the minimum  $MACR$ . This approach is simple, however it consumes additional network resources and may suffer from delay problems.

**Modifying the routers:** In this approach the source employs either Reno or Vegas [Jac88, BP95] and in addition it indicates its current rate (CR) in the IP (or TCP) header. Thus requiring (1) a modification of the TCP/IP header and (2) simple modifications in the routers. This modification enables the routers to identify sources whose current rate is relatively large (by comparing the sources current rate to the  $MACR$ ). We suggest four different ways by which the routers “tell” the source to decrease its rate in case of congestion:

1. **Explicit Rate Indication:** Simply implement in TCP a mechanism similar to the ABR rate base flow control. That is, the routers indicate the Explicit Rate (ER) on backward packets and the source station transmits according to that rate (in addition to the constraints imposed by the window size mechanism).
2. **Selective Discard:** The router discards any packet for which the indicated rate (CR) is larger than  $utilization\_factor \cdot MACR$ . This mechanism avoids congestion even in drop tail routers while reducing both the bias discussed in [FJ92] and the “beat down” problem. The right hand sides of Fig. 13 and Fig. 16 illustrate the behavior of this mechanism.

3. **Selective Source Quench:** A slight modification of the previous approach, where the router sends Source Quench messages to sources whose connection rate is above  $utilization\_factor \cdot MACR$ . The source reacts to the receipt of a Source Quench message as if a packet was dropped [BP87], and hence reduces its window size. Note that these messages might consume scarce network bandwidth at time of congestion. Alternatively, if an *EFCI* (Explicit Forward Congestion Indication) bit is used on the IP header, the router can set this bit. See [Flo94] for discussion of the use of a Source Quench mechanism and the addition of a new EFCI bit.
4. **Selective RED:** Here the router applies the RED mechanism (Random Early Detection). However, only packets whose rate is larger than  $utilization\_factor \cdot MACR$  may be dropped. Since we are selective in the packets we drop, we increase the fairness of the RED routers by reducing the “beat down” phenomena though this phenomena is not eliminated (in simulations, not presented here, it is seen that a session that crosses several routers is more vulnerable to such a drop even now).

If any of the suggested mechanisms is used in only part of the network routers then it is possible to operate both new and old methods in the network without that they prohibit each other. To this end, the “new” router contains two queues, *old-queue* for packets that do not follow the new algorithm and *new-queue* for packets that do. Packets are selected for transmission from the two queues in a round robin order.

Note, that the “new” TCP source code is an extension of either Vegas or Reno, and use their window adjustment policy. As a by product of being an extension, it correctly treat both old and new methods’ packets.

An additional motivation to implement the newly suggested flow control mechanism in TCP is that TCP traffic might traverse ATM networks. The use of a consistent flow control mechanism in both TCP and ABR over ATM, may improve the network utilization.

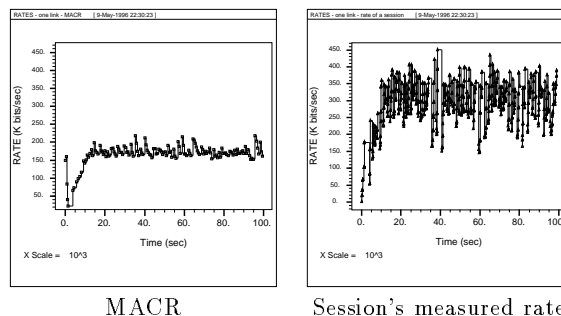


Figure 14: Four TCP sessions sharing a 1.5 Mbs link, each with a small RTT (40 ms) as in Fig. 13. We observe the values of  $MACR$  and the rate of an arbitrary session when the Selective discard mechanism is implemented. Note that the session rate is not stable due to the inherent burstiness of the window mechanism. However, the average rate is about twice the  $MACR$ .

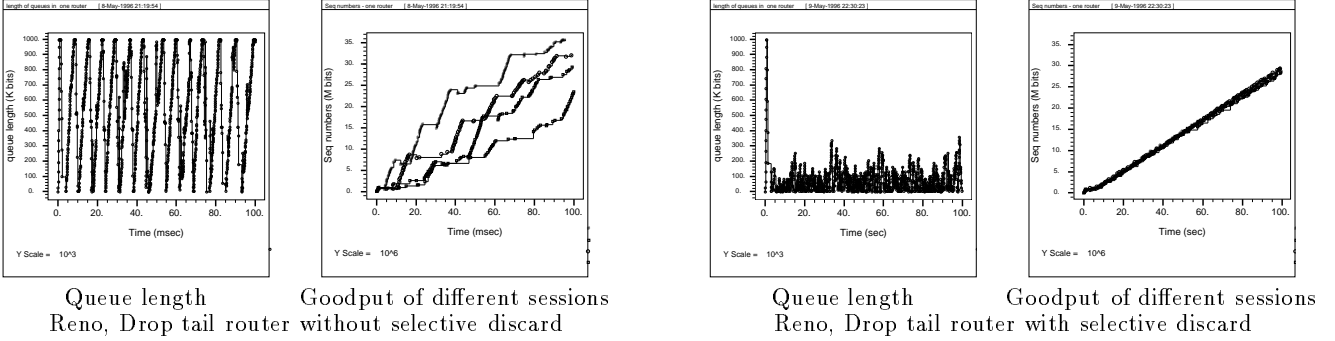


Figure 13: Four TCP sessions sharing a 1.5 Mbs link, each with a small RTT (40 ms). The two left hand diagrams illustrate the problematic case of a drop tail mechanism as is. Though all the sessions go via the same path, the link bandwidth is not equally shared (different sessions get different goodput). Note the sharp changes in the queue length. When the selective discard mechanism is used (two right hand side diagrams), the allocation is more fair and the queue length is moderate (400,000 bits, except for the initial peak).

### 3.3 Simulations of Selective Discard

We use BONES [ALT94] for our simulations. The TCP end systems implement Reno according to the pseudo code specified in Section 21 in [Ste94]. We assume greedy sources where size of packets is 512 bytes. Each source computes its rate as the ratio between the total amount of data received and acknowledged by the destination in a time interval, and the length of the time interval. The source computes that amount of data by tracking the difference in the number of bytes acknowledged. To avoid noisy peaks in this measurement we use the following standard exponential averaging:

**Parameters:**  $\tau$ ,  $\gamma$

**Initialization:**

$Rate := 0$ ;  $Curr\_Ack := 1$ ;  $MAX\_Ack := 1$

**Whenever an ACK is receipt:**

If  $ACK > Max\_Ack$  then  $Max\_Ack := ACK$

**At the end of each time interval of length  $\tau$ :**

$current\_rate := \frac{Max\_Ack - Curr\_Ack}{\tau}$   
 $Rate := Rate \cdot (1 - \gamma) + current\_rate \cdot \gamma$   
 $Curr\_Ack := Max\_Ack$

Where in our simulations  $\tau = 100ms$  and  $\gamma = 1/16$ .

The *goodput* of a session in any point of time is the total amount of “good” data that has been received by the sessions’ destination so far. We use the *goodput* measure for illustrations (Fig. 13 and Fig. 16) rather than the rate of traffic measured by the source at any point of time. This approach of illustration is used as the rate of TCP sessions is very bursty due to the manner in which the window is going up and down (see Fig. 14 for an illustration of the rate burstiness).

As for Phantom, the parameters are set as in Section 2 except for  $QT$  which is 8192 bits,  $\tau$  which is 100 msec and *ratio* (denoted also as  $f(\sigma)$ ) which is one quarter of the value in Section 2 (due to the window mechanism in TCP, the variance (burstiness) is much higher than in ABR). The *utilization\_factor* is set to 2, i.e., when using the selective discard mechanism, each packet with an indicated rate that is more than  $2 \cdot MACR$  is dropped.

Fig. 13 illustrates the behavior of the drop tail mechanism and the selective discard mechanism in a simple

configuration. Figure 14 shows the values of *MACR* and the rate of an arbitrary session when the selective discard mechanism is implemented.

Figure 15 presents a configuration where the “beat down” phenomena is likely to occur. Figure 16 illustrates the behavior of the drop tail mechanism and the selective discard mechanism in this configuration.

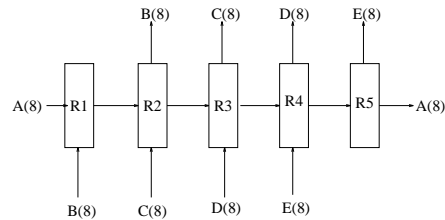


Figure 15: A configuration to measure the “beat down” phenomena. There are eight “long distance” sessions that cross four links and on each link there are eight additional “local” sessions. Each link has a bandwidth of 1.5 Mbs and a delay of 10 msec.

## 4 Comparison with other ATM rate based algorithms

We compare Phantom to four other constant space rate-based flow control algorithms that have been proposed in the ATM forum, EPRCA, suggested by Roberts [Rob94], APRC, suggested by Siu and Tzeng [ST94], CAPC, suggested by Barnhart [Bar94], and ERICA/ERICA+ suggested by Jain et. al. [JKVG95]. We give a brief description of each algorithm and compare it to Phantom. (The interested reader should look at a more complete survey by Jain [Jai95].)

All four algorithms compute for each output port in each switch a fair share parameter, (called *MACR* in EPRCA and APRC, *ERS* in CAPC, and *FS* in ERICA). This parameter is used to delimit the rate of sessions passing through the output port.

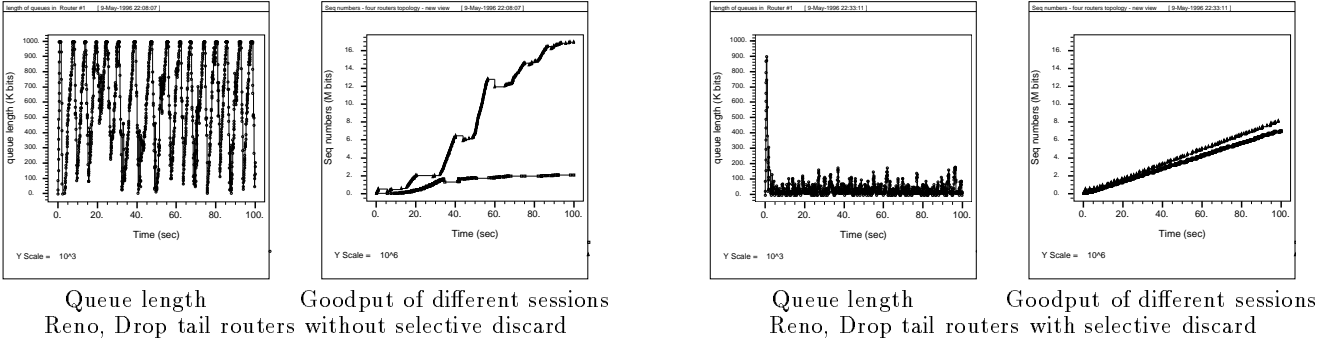


Figure 16: TCP sessions in the configuration described in Figure 15. In the right hand diagrams we see that the rate of a “local” session (group B) is much larger than the rate of a “long distance” session (group A). Note in the right hand side that the selective discard mechanism ensures a high level of fairness (almost equal allocation between “local” and “long distance” sessions) while keeping a moderate queue length.

#### 4.1 EPRCA and APRC

The EPRCA (Enhanced Proportional Rate Control Algorithm) [Rob94] was proposed by Roberts at the July 1994 ATM Forum meeting. In this algorithm the fair share parameter (MACR) is computed for each output port from the information received on the forward RM cells generated by the session’s source end system. The computation of MACR, and of the sessions’ rate limitations, depend on the queue state, which could be either *not congested*, if the queue length is smaller than a certain threshold, *congested*, if above the threshold, or *very congested*, if above a higher threshold. If the queue is *not congested*, the sessions’ rates are not limited. However, when the queue is either *congested* or *very congested* then the rates are restricted by some fraction of MACR. That is, the queue size plays a major role in the control loop. Furthermore, it influences the MACR indirectly through RM cells that go from the switch to the end-system and back to the switch. This extra delay in the control loop is a source for inherent oscillations in the algorithm as can be seen in Fig. 17. In the first part of Fig. 17 sessions increase their rate, when allowed, by a very small amount, while in the second part they may increase their rate more aggressively (the more aggressive increase rate is the one used in all other simulations in this paper). Aside from the oscillations, the extra delay in the control loop may cause an unfair sharing of the bandwidth between sessions with vastly different round trip delays [CGBS94, JKVG94, CRBdJ94].

In a modification of EPRCA, proposed by Siu and Tzeng called APRC (Adaptive Proportional Rate Control with intelligent congestion indication) [ST94] the definition of *congested* state is changed. Rather than being a function of the queue length it is now a function of the rate at which the queue length is changing. That is, if the derivative of the queue length is positive then the link is said to be *congested*.

Fig. 18 illustrates the behavior of APRC. Since the algorithm is insensitive to the absolute value of the queue length, the queue length might increase until it becomes *very congested* (because when in state *not congested* the queue length might increase more than it is decreasing while in state *congested*). Therefore, in some scenarios the queue length might often exceed the very congested threshold.

#### 4.2 CAPC

CAPC (Congestion Avoidance using Proportional Control) was proposed by Barnhart [Bar94]. It uses the fraction of unused capacity to control the algorithm actions. In this respect it is analogous to Phantom that uses the absolute amount of unused bandwidth to control the algorithm. CAPC estimate of the fair share is called ERS (Explicit Rate for Switch, which corresponds to our MACR). The main parameter used to control the changes of ERS is the ratio,  $r$ , between the rate of incoming traffic and the target link capacity. (Target link capacity is taken at 95% of the actual capacity.) Let  $\delta = 1 - r$ , if  $\delta$  is positive then ERS is multiplied by an increase factor, otherwise ERS is reduced by a multiplicative decrease factor. In addition, if the queue length exceeds a certain threshold, the switch instructs the source end-systems to reduce their rates (using the CI field, a binary indication provided by the ATM standard on RM cells).

Fig. 19 illustrates the behavior of CAPC in an environment with on/off sessions where the algorithm parameters are as recommended in [Bar94]. The configuration is analogous to that in Fig. 4 in Section 2. As can be seen CAPC has longer convergence time while its queue is relatively smaller during that time.

The larger value of the queue length in Phantom stems from the faster reaction of Phantom. We do not consider the momentary large queue as a problem, because Phantom is sensitive to the queue length (which enables it to smoothly and quickly drain the queue) and because the queue length consumed by Phantom can be reduced (as illustrated in Section 2).

Since CAPC uses the binary indication bit in very congested states, it is prone to the unfair behaviors indicated in [BdJ94]. Sessions that cross many switches are more likely to be “beaten down” and to get a smaller allocation than sessions that cross only a few switches. Unlike Phantom, in CAPC it is necessary to know the exact bandwidth allocated to ABR, a quantity which may change over time due to the addition and deletion of CBR and VBR sessions. A good feature that CAPC and some versions of Phantom share, is that they do not require forward RM cells.

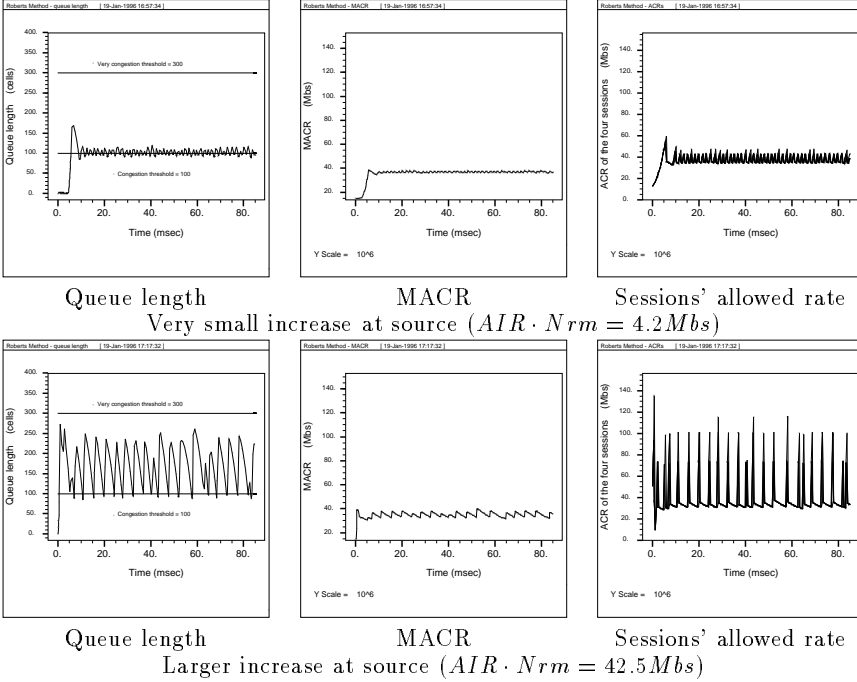


Figure 17: EPRCA - Four greedy sources sharing a 150Mbs link, RTT is negligible. Congested threshold is 100 cells and Very congested threshold is 300 cells, values of other parameters are as recommended in [Rob94].

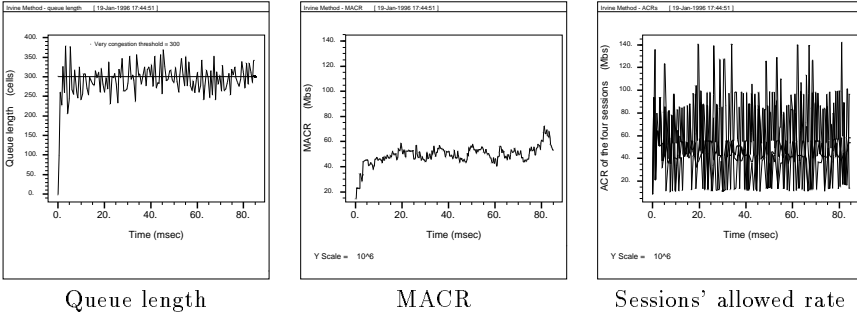


Figure 18: APRC - Four greedy sources sharing a 150Mbs link, RTT is negligible. Very congested threshold is 300 cells, values of other parameters are as recommended in [ST94].

### 4.3 ERICA and ERICA+

In [JKVG95, JKG<sup>+</sup>95] Jain et. al., from Ohio State University, presented the ERICA (Explicit Congestion Indication for Congestion Avoidance) and ERICA+ algorithms. In ERICA the switch computes two basic parameters: The *Load Factor*,  $z$  which equals to the input rate divided by the target rate (say 95% of the available bandwidth), and the *Fair Share*,  $FS$  which is the target rate divided by the number of active sessions. Given these two parameters the algorithm allows each session to change its rate to the maximum between  $FS$  and the session's current rate multiplied by  $\frac{1}{z}$ . Hence if the link is over loaded sessions drop their rate towards  $FS$  and if the link is under utilized sessions may increase their rates.

In ERICA+ the target rate is made a dynamic function of the queue length (actually of the delay data suffers in going through the queue). I.e., target rate = ratio · available bandwidth, such that ratio is decreased when

the queue length is increased and vice versa.

These schemes have several advantages. First, the resulting queue length is typically very small (in ERICA+), close to 1 and hence the delay is small too. Secondly, the algorithm reacts to changes in the traffic requirements relatively fast. Third, the algorithm is simple. Fig. 20 is an illustration of ERICA and ERICA+ where 24 sessions are sharing the same link (same settings as in Fig. 10 Section 2).

In [TW96] it is suggested that ERICA and ERICA+, may suffer from unfair bandwidth allocations (this example includes a network of more than a single link). I.e., even in a stable state the schemes might not necessarily reach the max-min fairness allocation. One possible reason for the unfairness is that  $FS$  as calculated by the algorithm is not the fair share that the link would have in a max-min fair allocation.

In order to achieve the fast convergence ERICA uses

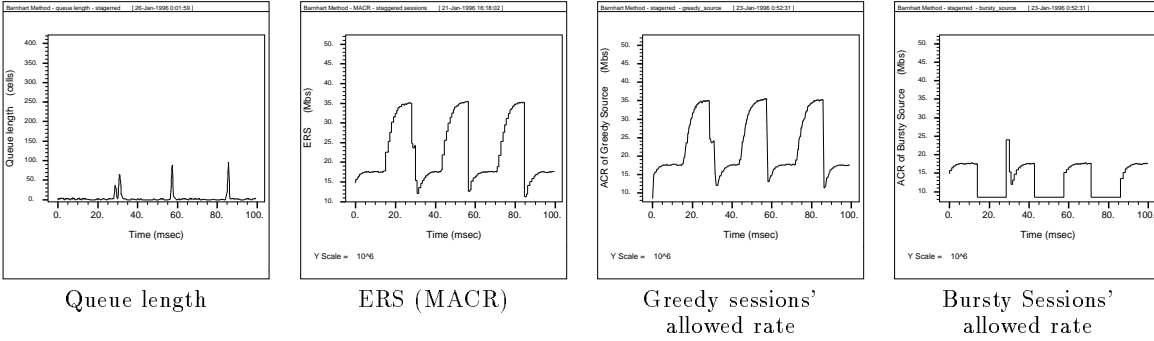


Figure 19: CAPC - Four greedy sessions and four on/off sessions sharing a 150 Mbps link with negligible RTT.

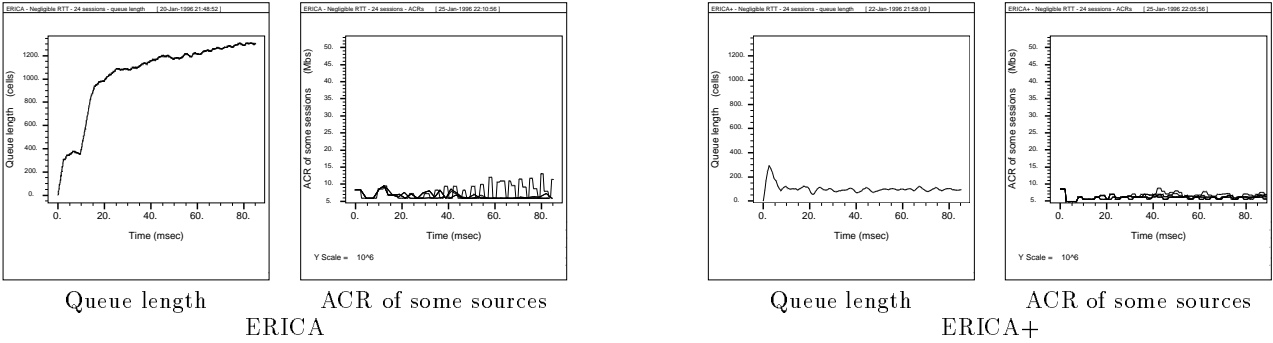


Figure 20: ERICA/ERICA+ - One link with 24 greedy sources where RTT is negligible.

a table of Current Cell Rates (CCRs) for all the sessions passing through a link, thus requiring a larger amount of space (non constant) and complicating the processing of RM cells. When implementing ERICA+ with constant space there is a scenario in which ERICA+ has an unfair allocation on a single link when there is a large difference in the delays of the sessions passing over the link [AMO96]. Finally, ERICA+ requires the explicit knowledge of the amount of available bandwidth.

### Acknowledgements

We would like to thank K.K. Ramakrishnan and the four anonymous referees for their very helpful comments.

### References

- [ALT94] ALTAGROUP. Bones designer core library reference. Technical report, December 1994.
- [AMO96] Y. Afek, Y. Mansour, and Z. Ostfeld. Phantom - a collection of simulation results. <ftp://ftp.math.tau.ac.il/pub/ostfeld/simulations/phantom.ps.Z>, May 1996.
- [Bar94] A. W. Barnhart. Explicit rate performance evaluation. Technical Report ATM-FORUM/94-0983R1, October 1994.
- [Bar95] A. W. Barnhart. Example switch algorithm for section 5.4 of tm spec. Technical Report ATM-FORUM/95-0195, Hughes Network Systems, February 1995.
- [BdJ94] J. C. R. Bennett and G. T. des Jardins. Comments on the july prca rate control baseline. Technical Report ATM-FORUM/94-0682, July 1994.
- [BF95] F. Bonomi and K. Fendick. The rate based flow control for available bit rate atm service. *IEEE Networks*, pages 24–39, March/April 1995.
- [BG87] D. P. Bertsekas and R. G. Gallager. *Data Networks*. Prentice Hall, 1987.
- [BP87] R. Braden and J. Postel. Requirements for internet gateways. Technical Report RFC 1009, IETF, June 1987.
- [BP95] L. S. Brakmo and L. L. Paterson. Tcp vegas: End to end congestion avoidance on a global internet. *IEEE JSAC*, 13,8:1465–1480, October 1995.
- [CCJ95] A. Charny, D.D. Clark, and R. Jain. Congestion control with explicit rate indication. In *Proc. ICC*, June 1995.
- [CGBS94] Y. Chang, N. Golmie, L. Benmohamed, and D. Su. Simulation study of the new rate-based eprca traffic management mechanism. Technical Report ATM-FORUM/94-0809, September 1994.
- [Cha94] A. Charny. An algorithm for rate allocation in a packet-switching network with feedback. Technical Report MIT/LCS/TR-601,

- MIT Laboratory for Computer Science, April 1994.
- [CR96] A. Charny and K.K. Ramakrishnan. Scalability issues for distributed explicit rate allocation in atm. In *Proc. 15th IEEE INFOCOMM*, pages 1182–1189, March 1996.
- [CRBdJ94] A. Charny, K.K. Ramakrishnan, J. C. R. Bennett, and G. T. des Jardins. Some preliminary results on the eprca rate-control scheme. Technical Report ATM-FORUM/95-0941, September 1994.
- [FJ92] S. Floyd and V. Jacobson. On traffic phase effects in packet-switches gateways. *Internetwork. Res. and Exper.*, 3(3):115–156, September 1992.
- [FJ93] S. Floyd and V. Jacobson. Random early detection gateways for congestion avoidance. *IEEE/ACM Transactions on Networking*, 1(4):397–413, August 1993.
- [Flo94] S. Floyd. Tcp and explicit congestion notification. *ACM Computer Communication Review*, 24,5:10–23, August 1994.
- [Gaf82] E. Gafni. *The Integration of Routing and Flow Control for Voice and Data in a Computer Communication Network*. PhD thesis, Massachusetts Institute of Technology, Dept. of Elec. Eng. and Comp. Science, Cambridge, MA, August 1982.
- [GB84] E. Gafni and D. P. Bertsekas. Dynamic control of session input rates in communication networks. *IEEE Trans. on Automatic Control*, AC-29(11):804–823, November 1984.
- [Hay81] H. Hayden. Voice flow control in integrated packet networks. Master’s thesis, Massachusetts Institute of Technology, Dept. of Electrical Engineering, Cambridge, MA, 1981.
- [Jac88] V. Jacobson. Congestion avoidance and control. In *Proc. SIGCOMM*, pages 314–329, August 1988.
- [Jaf81] J. M. Jaffe. Bottleneck flow control. *IEEE Trans on Communication*, COM-29,1(7):954–962, July 1981.
- [Jai95] R. Jain. Congestion control and traffic management in atm networks: Recent advances and a survey. *Computer Networks and ISDN Systems*, February 1995.
- [JKG<sup>+</sup>95] R. Jain, S. Kalyanaraman, R. Goyal, Sonia Fahmy, and Fang Lu. ERICA+: Extensions to the ERICA switch algorithm. Technical Report ATM-FORUM/95-1346, The Ohio State University, October 1995.
- [JKV94] R. Jain, S. Kalyanaraman, and R. Viswanathan. The OSU scheme for congestion avoidance using explicit rate indication. Technical Report ATM-FORUM/95-0883, The Ohio State University, September 1994.
- [JKVG94] R. Jain, S. Kalyanaraman, R. Viswanathan, and R. Goyal. Rate based schemes: Mistakes to avoid. Technical Report ATM-FORUM/94-0882, The Ohio State University, September 1994.
- [JKVG95] R. Jain, S. Kalyanaraman, R. Viswanathan, and R. Goyal. A sample switch algorithm. Technical Report ATM-FORUM/95-0178R1, The Ohio State University, February 1995.
- [Kes91] S. Keshav. A control-theoretic approach to flow control. In *Proc. SIGCOMM*, pages 3–16, September 1991.
- [KVR95] L. Kalampoukas, A. Varma, and K. K. Ramakrishnan. An efficient rate allocation algorithm for atm networks providing max-min fairness. In *Proc. HPN*, September 1995.
- [Mos84] J. Mosley. *Asynchronous Distributed Flow Control Algorithms*. PhD thesis, MIT, June 1984.
- [MR91] A. Mankin and K. K. Ramakrishnan. Gateway congestion control survey. Technical Report RFC 1254, IETF, August 1991.
- [RF94] A. Romanow and S. Floyd. Dynamics of tcp traffic over atm networks. In *Proc. SIGCOMM*, pages 79–88, August 1994.
- [RJ90] K. K. Ramakrishnan and R. Jain. A binary feedback scheme for congestion avoidance in computer networks. *ACM Trans. Comput System.*, 8(2):158–181, 1990.
- [Rob94] L. Roberts. Enhanced PRCA (proportional rate-control algorithm). Technical Report ATM-FORUM/94-0735R1, ATM Systems, August 1994.
- [Sat96] S. Sathaye. Atm forum traffic management specification version 4.0. Technical Report ATM-FORUM/95-0013R10, February 1996.
- [ST94] K. S. Siu and H. H. Tzeng. Adaptive proportional rate control (APRC) with intelligent congestion indication. Technical Report ATM-FORUM/94-0888, University of California, Irvine, September 1994.
- [Ste94] W. Stevens. *TCP/IP Illustrated*. Addison-Wesley, 1994.
- [TW96] D. H. K. Tsang and W. K. F. Wong. A new rate-based switch algorithm for ABR traffic to achieve max-min fairness with analytical approximation and delay adjustment. In *Proc. 15th IEEE INFOCOMM*, pages 1174–1181, March 1996.
- [Woj94] Lou Wojnaroski. Baseline text for traffic management sub-working group. Technical Report ATM-FORUM/94-0394R5, October 1994.
- [Zha89] L. Zhang. A new architecture for packet switching network protocols. Technical Report MIT/LCS/TR-455, MIT Laboratory for Computer Science, August 1989.