



Centrum voor Wiskunde en Informatica

State-dependent $M/G/1$ type queueing analysis for congestion control in data networks

E. Altman, K. Avrachenkov, C. Barakat, R. Núñez Queija

Probability, Networks and Algorithms (PNA)

PNA-R0005 July 31, 2000

Report PNA-R0005
ISSN 1386-3711

CWI
P.O. Box 94079
1090 GB Amsterdam
The Netherlands

CWI is the National Research Institute for Mathematics and Computer Science. CWI is part of the Stichting Mathematisch Centrum (SMC), the Dutch foundation for promotion of mathematics and computer science and their applications.

SMC is sponsored by the Netherlands Organization for Scientific Research (NWO). CWI is a member of ERCIM, the European Research Consortium for Informatics and Mathematics.

Copyright © Stichting Mathematisch Centrum
P.O. Box 94079, 1090 GB Amsterdam (NL)
Kruislaan 413, 1098 SJ Amsterdam (NL)
Telephone +31 20 592 9333
Telefax +31 20 592 4199

State-dependent M/G/1 Type Queueing Analysis for Congestion Control in Data Networks

Eitan Altman, Kostya Avrachenkov¹, Chadi Barakat²

INRIA

2004, route des Lucioles – B.P. 93, 06902 Sophia Antipolis, France

{Eitan.Altman,K.Avrachenkov,Chadi.Barakat}@sophia.inria.fr

Rudesindo Núñez Queija³

CWI

P.O. Box 94079, 1090 GB Amsterdam, The Netherlands

sindo@cwi.nl

ABSTRACT

We study a TCP-like linear-increase multiplicative-decrease flow control mechanism. We consider congestion signals that arrive in batches according to a Poisson process. We focus on the case when the transmission rate cannot exceed a certain maximum value. The distribution of the transmission rate in steady state as well as its moments are determined. Our model is particularly useful to study the behavior of TCP, the congestion control mechanism in the Internet. Burstiness of packet losses is captured by allowing congestion signals to arrive in batches. By a simple transformation, the problem can be reformulated in terms of an equivalent M/G/1 queue, where the transmission rate in the original model corresponds to the workload in the ‘dual’ queue. The service times in the queueing model are not i.i.d., and they depend on the workload in the system.

2000 Mathematics Subject Classification: 60K25, 68M20, 90B18, 90B22.

Keywords & Phrases: state-dependent queue, congestion control in data networks, linear increase and multiplicative decrease, transmission control protocol (TCP).

Note: The work of R. Núñez Queija is part of the PNA 2.1 project ‘Quality in Future Networks’ of the Telematics Institute.

¹ Supported with a CNET France-Telecom grant on flow control in High Speed Networks.

² Financed by an RNRT (French National Research Network in Telecommunications) “Constellations” project on satellite communications.

³ Visited INRIA Sophia Antipolis from February 1 until June 1 (2000) with support from the Netherlands Organization for Scientific Research (NWO).

1 Introduction

In today's high speed telecommunication networks, a large part of the traffic is able to adapt its rate to the congestion conditions in the network. Congestion control is typically designed so as to allow the transmission rate to increase linearly in time in the absence of congestion signals, whereas when congestion is detected, the rate decreases by a multiplicative factor. This is both the case of the Available Bit Rate (ABR) service category in ATM[1] (see definition and use of RDF and RIF) as well as the Transmission Control Protocol (TCP) in the Internet environment[9, 21]. Congestion is detected by the source through signals. In case of ABR, the congestion signals are RM (Resource Management) cells that have been marked due to congestion information in some switch along the path of the connection. In case of the Internet, the congestion signals are packet losses that are detected by the source either through the expiration of a retransmission timer, or through some negative acknowledgement mechanism (three duplicate ACKs [21]). There is also a proposal to add some explicit congestion signaling to the Internet (the ECN proposal [6]).

The performance evaluation of congestion control mechanisms is an important issue for network and protocol design. This evaluation requires a description of times between the arrivals of consecutive congestion signals. Experimentations over the Internet [4, 14] have shown that on long distance connections, the Poisson assumption about the times between congestion signals is quite reasonable. This happens when the throughput of the studied connection is small compared to the exogenous traffic, and when the number of hops on the path is large so that the superposition of the packet drops in routers leads to exponential times between congestion signals. For local area networks, we noticed that the congestion signals may arrive in bursts [4]. However, the times between bursts correspond well to the Poisson assumption. For this reason, we consider the case when congestion signals arrive in batches according to a Poisson process. Batches contain a random number of congestion signals and each such signal causes the division of the transmission rate by some constant γ . In the sequel, we also refer to a batch of congestion signals as a *loss event*.

We focus on the case when a certain limitation on the transmission rate exists. We determine the exact expression of the throughput under such a limitation. In the literature, only simplistic approximations have been proposed [3, 18] so far. We study two possible scenarios that lead to such a limitation:

(i) **Peak Rate limitation:** the limitation is not due to congestion in the network but rather to some external agreement. In that case, when the transmission rate reaches a certain level M , it remains constant until a loss event appears. For example in case of TCP, the window cannot exceed the buffer space available at the receiver [21]. In the ABR service of ATM, the transmission rate cannot exceed the Peak Cell Rate imposed by the contract between the user and the network. It is expected that such limitations on the transmission rate will become more and more important as the capacity and the speed of the links in the network grow, since it is then more likely that connections reach their maximum peak rate before congestion in the network occurs. Of course, this is not the case if peak rates increase in proportion with the speed of network links.

(ii) **Congestion limitation:** the limitation on the transmission rate is due to congestion in the network that occurs whenever the input rate reaches a level M . In that case we shall have an extra batch of congestion signals when the level M is attained which also causes a

decrease of the transmission rate by a random factor. A typical example of such limitation is the available bandwidth in the network. Another example is the reserved bandwidth in a Differentiated Services network [8] in cases where packets exceeding the reserved bandwidth are dropped rather than injected into the network as low priority packets [20].

In the particular case in which the batches contain a single congestion signal, the peak rate limitation model reduces to the one studied in [14], who already attempted at computing the first two moments of the transmission rate. A remarkable observation is done in that reference showing that the flow control can be reformulated in terms of an equivalent M/G/1 queue, where the transmission rate is translated into the workload of the queue. The congestion signals correspond to customers arriving at the queue according to a Poisson process. The service times in the ‘dual’ queueing model are not i.i.d., and they depend on the workload in the system. This transformation is also valid in our more general setting, except that in our model with congestion limitation, there is an additional arrival in the equivalent queueing model (in addition to the Poisson arrival stream) that occurs whenever the queue empties.

We solve the Kolmogorov equations and obtain the exact probability distribution as well as the moments of the transmission rate (of the window in case of TCP) for both problems. In doing so, we correct an error¹ in [14].

We briefly mention some related results. Queueing analysis with service times that depend on the workload or on the queue length have been also considered in [2, 11, 16, 19]. Our model is a special case of the one studied in [19], where an implicit characterization of the steady state distribution is obtained (closed-form expressions were obtained for special cases that do not cover our model). In [11] an asymptotic approximation is used for solving state-dependent GI/G/1 queues in which both inter-arrival times, service requirements and the service rate may depend on the workload. The peak rate limitation model is a special case of the model with a general stationary and ergodic arrival process studied in [4]. For that model only bounds on the throughput were obtained. Exact expressions for the throughput were obtained there for the case in which no limitation on the transmission rate exists (see also [3, 13, 15, 18]).

The paper is structured as follows. In Section 2 we describe a general model of flow control with limitation on the transmission rate and we provide a preliminary analysis. The two cases of peak rate limitation and congestion limitation are described separately in Sections 2.1 and 2.2. It is shown that a special case of the model with congestion limitation reduces to that of the model with peak rate limitation. Therefore, in the following we first focus on the case of peak rate limitation. In Section 3 we show that the model is dual to an M/G/1 queueing model with services that depend on the total workload in the system. We then derive the moments and the distribution of the transmission rate in Sections 4 and 5 in terms of the probability that the transmission rate is at its maximum value. This quantity can be determined using that the distribution function is non negative, but in order to derive a computationally tractable expression for it, we pursue an alternative approach in Section 6. The results are specified in Section 7 for an important particular case, that of one congestion signal per batch and a reduction factor of 2. This case corresponds to long distance TCP connections in today’s Internet where the congestion signals do not cluster significantly. In Section 8 we present the analysis for a more general model of congestion limitation than the

¹In a private communication, the authors of [14] announced to replace the draft.

one in Section 2.2. In the general case the model does not reduce to the peak rate limitation model. The model with peak rate limitation is validated in Section 9. By means of numerical examples, Section 9.1 illustrates that our results lend themselves for computation of the window (or, transmission rate) distribution and density functions. In Section 9.2 we compare the results of Section 7 to measurements from long distance TCP connections. Section 10 draws conclusions from the obtained results and indicates directions of further research. Finally, the Appendix displays technical results needed in the mathematical analysis.

2 Flow control with rate limitation: models and preliminary analysis

In this section we present our model for the rate evolution of the flow control mechanism. In the sequel we adopt the usual terminology for TCP, the well known window-based congestion control protocol of the Internet: we shall work with the window size rather than the transmission rate. The transmission rate of a window-based flow control mechanism is at any moment equal to the window size divided by the round-trip time (RTT) of the connection.

Let M denote the maximum window size. The limitation on the window size is either due to a peak rate limitation or to a congestion limitation. In the following we explain the similarities and the differences between the models in the two cases. While no congestion signal is received and the window is smaller than M , the window of the protocol increases linearly at rate $\alpha > 0$. In case of TCP, $\alpha = 1/(b \cdot RTT)$ where b is the number of data packets covered by an ACK (usually 2, see e.g.[18, 21]).

We assume that batches containing a random number of congestion signals arrive according to an independent Poisson process. We denote the sizes (i.e., the numbers of congestion signals) of consecutive batches by N_1, N_2, N_3, \dots , and we assume that these constitute an i.i.d. sequence. The size of an arbitrary batch is generically denoted by $N \stackrel{d}{=} N_k$. The Poisson process and the sequence $N_k, k = 1, 2, \dots$, are independent of each other and independent of the past evolution of the window. For each congestion signal received, the window is divided by a factor $\gamma > 1$ which is a fixed parameter. That is, if an arriving batch contains $N = n$ congestion signals, the window is multiplicatively decreased by a factor γ^{-n} . Immediately after the multiplicative decrease, the window restarts its linear increase. In case of peak rate limitation, the window stays constant at M when this maximum level is reached until the next congestion signal is received. In case of congestion limitation, immediately upon reaching M , a congestion signal is received and the window is decreased. We present the two cases separately in Sections 2.1 and 2.2, showing how the analysis of a particular case of the congestion limitation model reduces to that of the peak rate limitation model. In Section 8 we consider a more general model of congestion limitation.

First we introduce some further common notation. We denote the p.g.f. (probability generating function) of the distribution of N by

$$Q(z) := \mathbf{E} [z^N] =: \sum_{n=1}^{\infty} z^n q_n, \quad |z| \leq 1. \quad (2.1)$$

Note that the peak rate limitation model with $\gamma = 2$ and $q_1 = 1$ reduces to the model studied in [14], where congestion signals appear according to a Poisson process and where the window

is divided by two upon every congestion signal occurrence. By considering a general model, we aim to account for a wide range of flow control mechanisms other than TCP and for future enhancements to TCP congestion control.

Let us denote the window size at time $t \geq 0$ by $W(t) \in (0, M]$. We have the following stability result which follows from Theorem 1 in [4]:

Theorem 2.1 *There exists a stationary process $W^*(t)$ such that $W(t)$ converges to $W^*(t)$ in distribution for any initial state. Moreover, we have P -a.s.*

$$\lim_{t \rightarrow \infty} \sup_{s \geq t} |W(s) - W^*(s)| = 0. \quad (2.2)$$

Note that (2.2) implies that the stationary distribution of $W(t)$ is unique. For $x \in (0, M]$, denote the (time-average) distribution function by

$$F(x) := \lim_{T \rightarrow \infty} \frac{1}{T} \int_{t=0}^T \mathbf{P} \{W(t) \leq x\} dt. \quad (2.3)$$

It follows from Theorem 2.1 that this limit is independent of $W(0)$ and coincides with the stationary distribution of $W(t)$.

We first assume that $F(x)$ is continuous in $x \in (0, M)$ (in the case of peak rate limitation it is clear from physical considerations that $F(x)$ has an atom at $x = M$). Under this assumption we find a function $F(x)$ which is an equilibrium distribution for the window size and, hence, from its uniqueness it follows that it is the desired distribution. Instead of $F(x)$ it will be convenient to work with the complementary distribution function

$$\overline{F}(x) = 1 - F(x) = \mathbf{P} \{W > x\}, \quad x \in (0, M].$$

To differentiate between the cases of peak rate limitation and congestion limitation, in the latter case we attach a superscript ^{cl} to the symbols introduced above, e.g., the distribution function is denoted by $F^{\text{cl}}(x)$. Next we treat the two cases separately. We show how the analysis of a special case of the model with congestion limitation reduces to that of the model with peak rate limitation. The analysis of the general congestion limitation model is presented in Section 8.

2.1 Flow control with peak rate limitation

With peak rate limitation, when the window reaches the maximum level M , it stays there until the next congestion signal is received. In Section 3 below we show that the window size process $W(t)$ can be related to the workload of an M/G/1 queue (see also [14]). The workload of this state-dependent M/G/1 queue can be seen to be a Markov process (e.g., see [10]), and hence the window size evolution $W(t)$ is a Markov process as well. With this in mind, we derive a steady-state Kolmogorov equation for $\overline{F}(x) = \mathbf{P}\{W > x\}$ which is the basis to our analysis. We use the following up and down crossing argument: Assume that the process is in equilibrium and consider a level $x \in (0, M)$. Whenever the window size increases from less than or equal to x to more than x we say that an up crossing of the level x has occurred. Similarly, if the window size decreases from more than x to less than or equal

to x we say that a down crossing of the level x has occurred. Let $[t, t + \Delta]$ be a small time interval, where t is a deterministic time moment. When the process is in equilibrium, the probability of up-crossing

$$(1 - \lambda\Delta)\mathbf{P}\{x - \alpha\Delta < W \leq x\} + o(\Delta)$$

is equal to the probability of down-crossing

$$\lambda\Delta \sum_{n=1}^{\infty} q_n \mathbf{P}\{x < W \leq \min(\gamma^n x, M)\} + o(\Delta).$$

After equating these, we pass $\Delta \downarrow 0$. Since we assumed that $F(x) = \mathbf{P}\{W \leq x\}$ is continuous for $x < M$ (see Remark 5.1 for a justification of this assumption), we conclude that the derivative of $F(x)$ exists and is continuous for all x except at $x = M\gamma^{-n}$, when $q_n > 0$. For $x \in (0, M) \setminus \{M\gamma^{-n}\}_{n=1,2,\dots}$ we obtain the following steady-state Kolmogorov equation

$$\alpha \frac{d}{dx} \mathbf{P}\{W \leq x\} = \lambda \sum_{n=1}^{\infty} q_n \mathbf{P}\{x < W \leq \min(\gamma^n x, M)\},$$

or, equivalently,

$$-\alpha \frac{d}{dx} \bar{F}(x) = \lambda \left(\bar{F}(x) - \sum_{n=1}^{\infty} q_n \bar{F}(\min(\gamma^n x, M)) \right). \quad (2.4)$$

From this differential equation we shall determine $\bar{F}(x)$, $x \in (0, M)$, in terms of the probability

$$P_M := \mathbf{P}\{W = M\} = 1 - F(M-) = \bar{F}(M-).$$

In Section 4 we first use (2.4) to determine the moments of the window size distribution in terms of P_M . Then we find the distribution function itself in Section 5. The unknown P_M is then determined using the fact that $\bar{F}(x)$ is a complementary probability distribution function ($\bar{F}(0) = 1$). However, the expression obtained for P_M in this way, does not lend itself for computational purposes. Therefore we show an elegant alternative to determine P_M in Section 6, which leads to an efficient and numerically stable algorithm for computations.

2.2 Congestion limitation: a special case

When the maximum window size M is due to congestion limitation, immediately upon reaching the level M a batch of congestion signals is generated. In this section we study the case when the size of such a batch has the same distribution as the random variable N . In Section 8 we present the analysis of the more general case when the number of congestion signals that result from reaching M has a different distribution than N . Similarly as in Section 2.1, we can derive the following differential equation for $\bar{F}^{\text{cl}}(x)$, $0 < x < M$:

$$\begin{aligned} -\alpha \frac{d}{dx} \bar{F}^{\text{cl}}(x) = & \lambda \left(\bar{F}^{\text{cl}}(x) - \sum_{n=1}^{\infty} q_n \bar{F}^{\text{cl}}(\min(\gamma^n x, M)) \right) \\ & + \lambda g \mathbf{P} \left\{ N \geq \frac{\ln(M) - \ln(x)}{\ln(\gamma)} \right\}, \end{aligned} \quad (2.5)$$

with,

$$g := -\frac{\alpha}{\lambda} \frac{d}{dy} \overline{F}^{\text{cl}}(y) \Big|_{y=M-}.$$

The additional term, compared to (2.4), comes from the fact that a down crossing of the level x may be due to the fact that the level M is reached and that the rate is decreased by a factor γ^{-n} with $\gamma^{-n}M \leq x$. Note that if $\overline{F}(x)$ is the unique complementary distribution function satisfying (2.4) then

$$\overline{F}^{\text{cl}}(x) := \frac{\overline{F}(x) - P_M}{1 - P_M}, \quad 0 < x < M, \quad (2.6)$$

is the unique complementary distribution function satisfying (2.5). This follows immediately by substituting (2.6) into (2.5). This relation has a simple geometric interpretation. Using the fact that the Poisson process is memoryless, if we consider the model with peak rate limitation only at moments when the window is less than M (i.e., we cut out all periods where the window equals M), what we get is identical to the model with congestion limitation. Thus, we can concentrate on finding the distribution function $F(x)$ for the peak rate limitation model and then use (2.6) or the equivalent:

$$F^{\text{cl}}(x) = \frac{F(x)}{1 - P_M}.$$

In particular, the moments of the window size in the two models are related by:

$$\mathbf{E} \left[\left(W^{\text{cl}} \right)^k \right] = \frac{\mathbf{E} [W^k] - P_M M^k}{1 - P_M}. \quad (2.7)$$

In Section 4 below we derive a recursive relation for $\mathbf{E} [W^k]$. Combined with (2.7), this gives a recursion on $\mathbf{E} \left[\left(W^{\text{cl}} \right)^k \right]$ which we report at this point for completeness:

$$\mathbf{E} \left[\left(W^{\text{cl}} \right)^k \right] = \frac{k\alpha \mathbf{E} \left[\left(W^{\text{cl}} \right)^{k-1} \right]}{\lambda (1 - Q(\gamma^{-k}))} - \frac{P_M}{1 - P_M} M^k. \quad (2.8)$$

Remark 2.1 We emphasize that in the congestion limitation model, the quantity P_M has no clear interpretation. In Section 6 we use the interpretation of this quantity in the peak rate limitation model to compute it. If we were to analyze the congestion limitation model without using (2.6), then from (2.5) we could express $\overline{F}^{\text{cl}}(x)$ — using the same techniques as in Section 5 — in terms of g instead of P_M . Note that these two constants are related:

$$g = \frac{P_M}{1 - P_M}, \quad (2.9)$$

The constant g can be determined using that $\overline{F}^{\text{cl}}(x)$ is a complementary probability distribution, see (5.8) below. Since from the analysis of Section 6 we obtain a more tractable expression for P_M (see Remark 5.2 for a related discussion), we will not further dwell on this approach.

3 The dual queueing model

Before proceeding with determining the moments and the distribution of the window size, we briefly show how the problem can be related to an M/G/1 queueing problem with service depending on the system workload, see also [14]. First we concentrate on peak rate limitation, below we comment on congestion limitation. Define

$$U(t) = \frac{M - W(t)}{\alpha}. \quad (3.1)$$

I.e., $U(t)$ is obtained by ‘flipping’ $W(t)$ around a horizontal line and then scaling by a factor $1/\alpha$. In particular, the area between $W(t)$ and the maximum window size M (Figure 1) corresponds to the area below $U(t)$. Note that $U(t)$ resembles the evolution in time of the workload (or the virtual waiting time) in a queueing system. A window equal to M corresponds to an empty queueing system. The linear increase in workload between arrivals of congestion signals corresponds to the decrease in workload due to service in the queueing model. The arrival of a batch of congestion signals in our model corresponds to an arrival to the queue. The reduction of the window upon a loss event corresponds to the increase in workload upon arrival in the equivalent queueing model. Given that the amount by which the window is reduced depends on the current value of the window (and of course on the number of congestion signals in the batch), the service time in the dual queueing model is dependent on the current workload there. We conclude that the dual model behaves indeed as an M/G/1 queue (infinite buffer capacity, one server and Poisson arrivals with intensity λ) with state-dependent service requirements. If U_n is the workload seen by arrival n in the M/G/1 queue, then its service time x_n is equal to

$$x_n = \left(\frac{M}{\alpha} - U_n \right) \cdot \left(1 - \frac{1}{\gamma^{N_n}} \right),$$

where N_n is the number of congestion signals in the n th batch of congestion signals in the original model. Instead of directly working with the congestion control model as we do in this paper, one could analyze the queueing model and switch back to the flow control problem by using Equation (3.1). In particular, $\mathbf{E}[W^k] = \mathbf{E}[(M - \alpha U)^k]$, $\mathbf{P}\{W \leq x\} = 1 - \mathbf{P}\{U \leq (M - x)/\alpha\}$ for $x < M$ and P_M is equal to the fraction of time that the dual queue is empty.

In the case of congestion limitation, the only difference in the dual queueing model is that we have an additional arrival once the system becomes empty. Thus, the arrival process is the sum of a Poisson process of intensity λ and another process that depends on the workload of the system (it generates an immediate arrival when the queue becomes empty). The definition of the service times in the dual queue and the transformation back to the flow control problem remain the same.

4 Moments of the window size distribution

Now focus on the model with peak rate limitation (the results obtained can also directly be used for the special case of congestion limitation described in Section 2.2). In this section we study the moments of the window size. The k -th moment of the transmission rate can be

simply obtained by dividing the k -th moment of the window size by $(RTT)^k$. Of particular interest is the expectation of the transmission rate which coincides with the throughput of the transfer or the time average of the transmission rate. Let \bar{X} denote the throughput. We have

$$\bar{X} = \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T X(t) dt = \frac{\mathbf{E}[W]}{RTT}. \quad (4.1)$$

Define for $\text{Re}(\omega) \geq 0$ the LST (Laplace-Stieltjes Transform) of the window size distribution by

$$\hat{f}(\omega) = \int_{x=0}^{M+} e^{-\omega x} dF(x).$$

Taking LTs (Laplace Transforms) in (2.4) leads to:

$$\alpha \left(\hat{f}(\omega) - P_M e^{-\omega M} \right) = \lambda \frac{1 - \hat{f}(\omega)}{\omega} - \lambda \sum_{n=1}^{\infty} \gamma^{-n} q_n \frac{1 - \hat{f}(\gamma^{-n} \omega)}{\gamma^{-n} \omega}. \quad (4.2)$$

Note that (4.2) holds in particular for $M = \infty$, i.e., no limitation on the window size, in which case $P_M = 0$. Using $\mathbf{E}[W^k] \leq M^k$, $k = 1, 2, \dots$, we may write

$$\begin{aligned} \hat{f}(\omega) &= 1 + \sum_{k=1}^{\infty} \frac{(-\omega)^k}{k!} \mathbf{E}[W^k], \\ \frac{1 - \hat{f}(\gamma^{-n} \omega)}{\gamma^{-n} \omega} &= \sum_{k=0}^{\infty} \frac{(-\gamma^{-n} \omega)^k}{(k+1)!} \mathbf{E}[W^{k+1}]. \end{aligned}$$

Substituting this in (4.2), using the absolute convergence of the doubly-infinite series to interchange the order of summation and equating the coefficients of equal powers of ω we get, for $k = 1, 2, \dots$,

$$\mathbf{E}[W^k] = \frac{k\alpha (\mathbf{E}[W^{k-1}] - P_M M^{k-1})}{\lambda(1 - Q(\gamma^{-k}))}, \quad (4.3)$$

from which the moments of the window size distribution can be obtained recursively. In particular we find for $k = 1, 2$:

$$\mathbf{E}[W] = \frac{\alpha(1 - P_M)}{\lambda(1 - Q(\gamma^{-1}))}, \quad (4.4)$$

$$\mathbf{E}[W^2] = \frac{2\alpha [\alpha(1 - P_M) - \lambda P_M M (1 - Q(\gamma^{-1}))]}{\lambda^2 (1 - Q(\gamma^{-1})) (1 - Q(\gamma^{-2}))}. \quad (4.5)$$

These first two moments can also be obtained using direct arguments, see Remarks 4.1 and 4.2 below. Such arguments were also used by Misra et al. [14] for the case $\gamma = 2$ and $N \equiv 1$. However, in their analysis an error appears which results in an additional equation besides

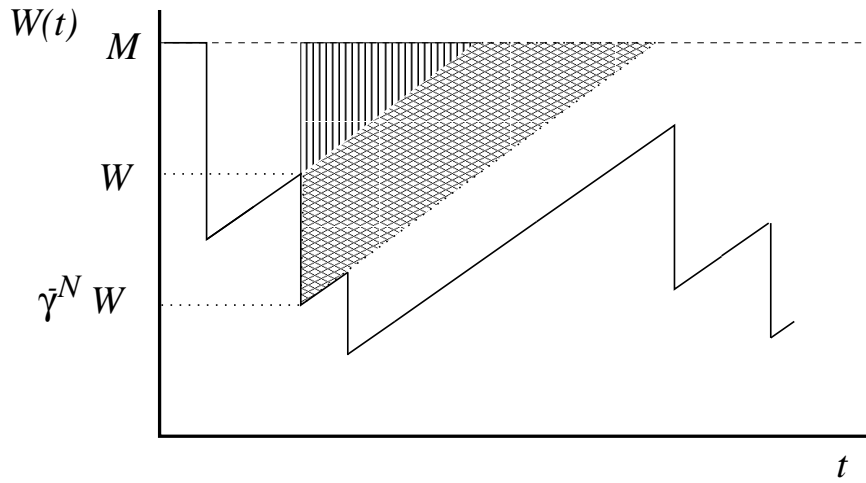


Figure 1: Area associated with a single loss

(4.4) and (4.5) from which they determine an incorrect expression for the probability P_M (see Remark 4.2).

Remark 4.1 The mean window size can be obtained by considering the mean drift. The upward drift of the window size is given by $\alpha \mathbf{P}\{W < M\}$ and the downward drift equals $\lambda \mathbf{E}[W] (1 - \mathbf{E}[\gamma^{-N}])$. Equating these gives (4.4).

We can further derive $\mathbf{E}[W^2]$ applying an argument similar to Little's law as was done by Misra et al. [14] for the case $\gamma = 2$ and $N \equiv 1$. The main idea is sketched in the following. For the dual queueing model described in Section 3, we can equate the mean workload $\mathbf{E}[U]$ with λ times the mean area below $U(t)$ 'induced by a single arrival' (use that Poisson arrivals see time averages: PASTA). Back in the original model, the 'mean surface' of the area *above* $W(t)$ in Figure 1 equals $M - \mathbf{E}[W]$. We find an alternative expression for this area by determining the surface of the area 'induced' by a loss event. This is also depicted in Figure 1. Suppose a loss occurs at window size W and the window is reduced by a factor γ^{-N} . We can associate with this loss an area above the curve (the surface of the larger triangle minus that of the smaller one) equal to

$$\frac{1}{2\alpha} (M - \gamma^{-N} W)^2 - \frac{1}{2\alpha} (M - W)^2.$$

Because of PASTA and the fact that N is independent of W , the expectation of the surface of this area is

$$\frac{1}{2\alpha} ((Q(\gamma^{-2}) - 1) \mathbf{E}[W^2] - 2M (Q(\gamma^{-1}) - 1) \mathbf{E}[W]).$$

The rate at which losses occur is λ , and so:

$$M - \mathbf{E}[W] = \frac{\lambda}{2\alpha} ((Q(\gamma^{-2}) - 1) \mathbf{E}[W^2] - 2M (Q(\gamma^{-1}) - 1) \mathbf{E}[W]). \quad (4.6)$$

Together with (4.4) this indeed gives (4.5).

Remark 4.2 For a special case of our model, yet another way is pursued in [14] to derive (4.4) and (4.5). However, there, the final result is incorrect due to a small error in an intermediate step. Defining $P_M(t) := \mathbf{P}\{W(t) = M\}$, the first two moments of $W(t)$ satisfy:

$$\begin{aligned} \frac{d}{dt} \mathbf{E}[W(t)] &= -\lambda(1 - Q(\gamma^{-1})) \mathbf{E}[W(t)] \\ &\quad + \alpha(1 - P_M(t)), \\ \frac{d}{dt} \mathbf{E}[W(t)^2] &= -\lambda(1 - Q(\gamma^{-2})) \mathbf{E}[W(t)^2] \\ &\quad + 2\alpha(\mathbf{E}[W(t)] - MP_M(t)). \end{aligned}$$

In steady state we have $\mathbf{E}[W(t)] \equiv \mathbf{E}[W]$, $\mathbf{E}[W(t)^2] \equiv \mathbf{E}[W^2]$ and $P_M(t) \equiv P_M$. Substitution into (4.7) gives (4.4) and substitution into (4.7) gives:

$$0 = -\lambda(1 - Q(\gamma^{-2})) \mathbf{E}[W^2] + 2\alpha(\mathbf{E}[W] - P_M M). \quad (4.7)$$

The latter is a linear combination of (4.4) and (4.6) and, hence, leads to (4.5). For the case $\gamma = 2$ and $N \equiv 1$, the formula given in [14] for $\mathbf{E}[W^2]$ (below Formula (4) in that reference) differs from (4.7) by a factor $-\alpha = -1/RTT$. This resulted in a third (incorrect) equation which is independent of (4.4) and (4.5) from which P_M was determined simultaneously with $\mathbf{E}[W]$ and $\mathbf{E}[W^2]$. In Section 6 we show how P_M can be determined correctly and computed efficiently.

5 Window size distribution function

In this section we determine the cumulative distribution function of the window size explicitly. The distribution of the transmission rate can be simply obtained by rescaling the window axis by $1/RTT$. We start with the case where $M < \infty$, providing an expression of the distribution function in the intervals $[M/\gamma^k, M/\gamma^{k-1}]$ with $k = 1, 2, \dots$. Then, for the case $M = \infty$, we give an expression of the distribution for any $x > 0$ as an infinite sum of exponentials.

5.1 Window distribution for finite M

For $M/\gamma \leq x < M$, Equation (2.4) reduces to:

$$-\alpha \frac{d}{dx} \bar{F}(x) = \lambda \bar{F}(x),$$

hence,

$$\bar{F}(x) = P_M e^{\frac{\lambda}{\alpha}(M-x)}, \quad \frac{M}{\gamma} \leq x < M. \quad (5.1)$$

To find the entire distribution we introduce, for $k = 1, 2, 3, \dots$,

$$\bar{F}_k(x) := \bar{F}(x), \quad \frac{M}{\gamma^k} \leq x < \frac{M}{\gamma^{k-1}}. \quad (5.2)$$

Equation (2.4) can now be written as:

$$\frac{d}{dx} \overline{F}_k(x) = -\frac{\lambda}{\alpha} \overline{F}_k(x) + \frac{\lambda}{\alpha} \sum_{n=1}^{k-1} q_n \overline{F}_{k-n}(\gamma^n x). \quad (5.3)$$

Since $\overline{F}(x)$ is continuous for $0 < x < M$ we have:

$$\overline{F}_k\left(\frac{M}{\gamma^{k-1}}\right) = \overline{F}_{k-1}\left(\frac{M}{\gamma^{k-1}}\right), \quad k = 2, 3, \dots \quad (5.4)$$

\overline{F}_k is recursively given by

$$\begin{aligned} \overline{F}_k(x) &= \overline{F}_{k-1}(M/\gamma^{k-1}) e^{\frac{\lambda}{\alpha} \left(\frac{M}{\gamma^{k-1}} - x\right)} \\ &\quad - \frac{\lambda}{\alpha} e^{-\frac{\lambda}{\alpha} x} \int_{u=x}^{M/\gamma^{k-1}} e^{\frac{\lambda}{\alpha} u} \sum_{n=1}^{k-1} q_n \overline{F}_{k-n}(\gamma^n u) du. \end{aligned}$$

We conclude from the above recursion that a solution to (5.3) and (5.4) has the following form

$$\overline{F}_k(x) = P_M \sum_{i=1}^k c_i^{(k)} e^{-\frac{\lambda}{\alpha} \gamma^{i-1} x}, \quad k = 1, 2, \dots \quad (5.5)$$

To determine the coefficients $c_i^{(k)}$, we substitute (5.5) into (5.3) and change the order of summation in the last term $\frac{\lambda}{\alpha} \sum_{n=1}^{k-1} q_n \overline{F}_{k-n}(\gamma^n x)$:

$$\begin{aligned} P_M \sum_{i=1}^k c_i^{(k)} \left(-\frac{\lambda}{\alpha} \gamma^{i-1}\right) e^{-\frac{\lambda}{\alpha} \gamma^{i-1} x} &= \\ &= -\frac{\lambda}{\alpha} P_M \sum_{i=1}^k c_i^{(k)} e^{-\frac{\lambda}{\alpha} \gamma^{i-1} x} + \frac{\lambda}{\alpha} P_M \sum_{i=1}^{k-1} \left[\sum_{n=1}^i q_n c_{i-n+1}^{(k-n)} \right] e^{-\frac{\lambda}{\alpha} \gamma^i x} \end{aligned}$$

By equating the terms with the same exponents, we get the following recursive formula

$$c_{i+1}^{(k)} = \frac{1}{1 - \gamma^i} \sum_{n=1}^i q_n c_{i-n+1}^{(k-n)}, \quad i = 1, \dots, k-1. \quad (5.6)$$

Once the coefficients $c_i^{(k)}$, $i = 2, \dots, k$ are computed, the coefficient $c_1^{(k)}$ can be determined from (5.4):

$$c_1^{(k)} = e^{\frac{\lambda}{\alpha} \frac{M}{\gamma^{k-1}}} \left[\sum_{i=1}^{k-1} c_i^{(k-1)} e^{-\frac{\lambda}{\alpha} \gamma^{i-k} M} - \sum_{i=2}^k c_i^{(k)} e^{-\frac{\lambda}{\alpha} \gamma^{i-k} M} \right]. \quad (5.7)$$

Note that to compute the coefficients $c_i^{(k)}$, we *do not* need P_M . Hence, using that $\overline{F}(x)$ is a complementary distribution function, P_M is then determined by:

$$1 = \overline{F}(0) = \lim_{k \rightarrow \infty} \overline{F}_k(M/\gamma^{k-1}) = P_M \left(\lim_{k \rightarrow \infty} \sum_{i=1}^k c_i^{(k)} e^{-\frac{\lambda}{\alpha} M/\gamma^{k-i}} \right). \quad (5.8)$$

However, this relation is not suitable to compute P_M , see Remark 5.2 below.

Remark 5.1 With (5.2) and (5.5) we have found an equilibrium distribution function $F(x)$ satisfying (2.4). By Thm. 2.1 it is the unique solution and, hence, the assumption that $F(x)$ is continuous for $x < M$ is justified.

Remark 5.2 Recursion (5.6) is suitable to determine the distribution function on an interval $M/\gamma^k \leq x \leq M$ when k is not too large. For large k the recursion may become instable, since it involves subtraction of numbers of the same order. Therefore (5.8) is not suitable to compute P_M . In Section 6 below we derive an alternative expression for P_M , which leads to a numerically stable and efficient algorithm to compute P_M .

Remark 5.3 One can alternatively show that the functions $\overline{F}_k(x)$ are of the form (5.5) using Laplace Transform techniques. Note that by means of the differential equations (5.3), these functions can be extended beyond the intervals $[M/\gamma^k, M/\gamma^{k-1}]$ to the whole real line. Of course, outside the intervals $[M/\gamma^k, M/\gamma^{k-1}]$ the functions $\overline{F}_k(x)$ may (and will) be unequal to \overline{F} . One may take Laplace Transforms in (5.3), solve the resulting recursion on k , and invert the transforms after applying partial fraction expansion. This approach is used in Sections 5.2 and 6.

5.2 Window distribution for infinite M

In this case, the results derived in the previous subsection cannot be applied immediately by letting M go to infinity. However, we can derive the LST of the window size distribution by similar arguments as before. When $M = \infty$, (4.2) becomes

$$\hat{f}(\omega) = -\frac{\lambda}{\alpha} \left[\frac{\hat{f}(\omega)}{\omega} - \sum_{n=1}^{\infty} q_n \frac{\hat{f}(\gamma^{-n}\omega)}{\omega} \right],$$

or, equivalently,

$$\hat{f}(\omega) = \frac{\frac{\lambda}{\alpha}}{\omega + \frac{\lambda}{\alpha}} \sum_{n=1}^{\infty} q_n \hat{f}(\gamma^{-n}\omega). \quad (5.9)$$

Substituting the above equation repeatedly into itself l times, applying partial fraction expansion at each step, and then taking $l \rightarrow \infty$, we conclude that $\hat{f}(\omega)$ can be expressed as follows:

$$\hat{f}(\omega) = \sum_{i=0}^{\infty} c_i \frac{-\frac{\lambda}{\alpha} \gamma^i}{\omega + \frac{\lambda}{\alpha} \gamma^i}, \quad (5.10)$$

for certain coefficients c_i (this is formally justified later). To determine the constants c_i , $i = 0, 1, \dots$, we substitute (5.10) into (5.9) and equate coefficients multiplying the terms $1/(\omega +$

$\frac{\lambda}{\alpha}\gamma^i$). This leads to the recursive formula

$$\frac{c_i}{c_0} = \frac{1}{1 - \gamma^i} \sum_{k=1}^i q_k \frac{c_{i-k}}{c_0},$$

which determines the ratios c_i/c_0 (it is for this reason that both sides contain a factor $1/c_0$). The coefficient c_0 follows from $f(0) = -\sum_{i=0}^{\infty} c_i = 1$:

$$c_0 = -\left(1 + \sum_{i=1}^{\infty} \frac{c_i}{c_0}\right)^{-1}. \quad (5.11)$$

Inversion of (5.10) back into the time domain gives:

$$F(x) = C + \sum_{i=0}^{\infty} c_i e^{-\frac{\lambda}{\alpha}\gamma^i x}, \quad (5.12)$$

with $C = 1$ (because $F(0) = 0$). Note that the above series is absolutely convergent for any value of $x \in [0, \infty)$. Thus, it is the (unique!) solution to (2.4) when $M = \infty$. For the case of no window size limitation and $N \equiv 1$, ($q_1 = 1$), $F(x)$ was already obtained in [17].

6 The probability of maximum window size

In Sections 4 and 5 we determined the window size distribution and its moments in terms of P_M . In this section we derive an expression for P_M from which it can be computed efficiently. For this we introduce the random variable $T(x)$ which is the time until the window size returns to the value x , starting just after a loss event occurs with the window size being equal to $x \in (0, M]$. We denote its expectation by $E(x) := \mathbf{E}[T(x)]$, $x \in (0, M]$. Then, from elementary renewal theory,

$$P_M := \mathbf{P}\{W = M\} = \frac{1/\lambda}{1/\lambda + E(M)}. \quad (6.1)$$

We now proceed to find the function $E(x)$. A typical evolution of the window size is depicted in Figure 2. For simplicity in the figure only losses having $N = 1$ are depicted and the times to recover from losses are partly cut out of the picture (denoted by the shaded areas). Suppose for the moment that the initial loss (at the level x) was such that $N = n$ (in the figure $n = 1$). Let $T_n(x)$ be the time to get back at level x conditional on $N = n$ and we further write $E_n(x) := \mathbf{E}[T_n(x)] := \mathbf{E}[T(x)|N = n]$. Note that

$$E(x) = \sum_{n=1}^{\infty} q_n E_n(x). \quad (6.2)$$

If no losses occur during the time $T_n(x)$ then $T_n(x) = (1 - \gamma^{-n})x/\alpha$, i.e., the window size x is reached in a straight line from the starting point at $\gamma^{-n}x$ (in the figure $\gamma^{-1}x$). Each time a loss occurs at a level $y \in (\gamma^{-n}x, x)$ it takes $T(y)$ time units to get back at the level y . Because of the memoryless property of the Poisson process, if we take out the shaded areas in Figure

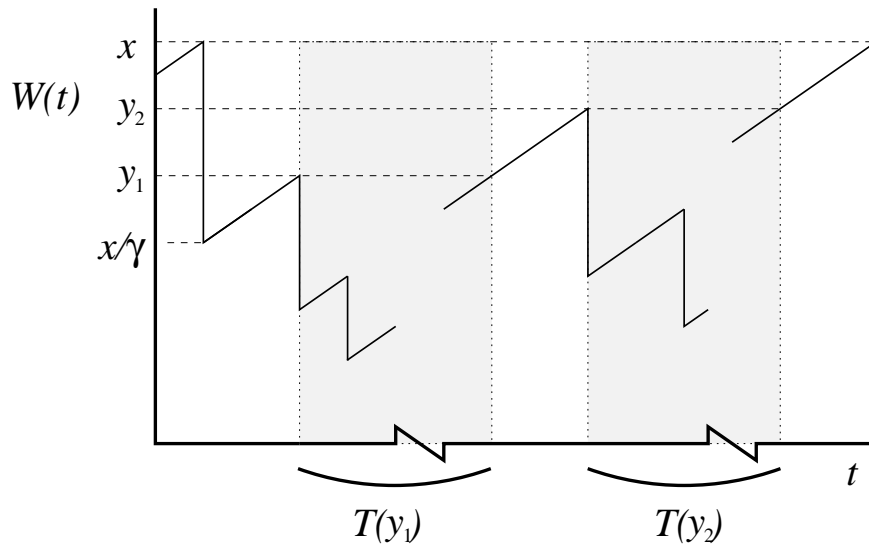


Figure 2: TCP window

2 and concatenate the non-shaded areas then the cut points (where the shaded areas used to be) form a Poisson process on the straight line from $\gamma^{-n}x$ to x . Thus if the cut points are given by y_1, y_2, \dots, y_m (in the figure $m = 2$) then

$$E_n(x) = \frac{(1 - \gamma^{-n})x}{\alpha} + E(y_1) + E(y_2) + \dots + E(y_m).$$

Since the loss process is a Poisson process, the mean number of cut points is $\lambda(1 - \gamma^{-n})x/\alpha$ and the position of each of the points y_j is uniformly distributed over the interval $(\gamma^{-n}x, x)$, see for instance [22, Thm. 1.2.5]. Hence,

$$\begin{aligned} E_n(x) &= \frac{(1 - \gamma^{-n})x}{\alpha} + \lambda \frac{(1 - \gamma^{-n})x}{\alpha} \int_{y=\gamma^{-n}x}^x \frac{E(y)}{(1 - \gamma^{-n})x} dy \\ &= \frac{(1 - \gamma^{-n})x}{\alpha} + \frac{\lambda}{\alpha} \int_{y=\gamma^{-n}x}^x E(y) dy. \end{aligned} \quad (6.3)$$

Using (2.1) and (6.2) we now arrive at

$$E(x) = \frac{(1 - Q(\gamma^{-1}))x}{\alpha} + \frac{\lambda}{\alpha} \sum_{n=1}^{\infty} q_n \int_{y=\gamma^{-n}x}^x E(y) dy. \quad (6.4)$$

Although in the finite-window case ($M < \infty$) the above integral equation has only meaning for $0 < x \leq M$, it is well defined for all $x > 0$. In the following we solve the integral equation for all $x > 0$. First we note that it has a unique solution, see Appendix A.1. Define the LT (Laplace Transform) of $E(x)$:

$$\hat{e}(\omega) := \int_{x=0}^{\infty} e^{-\omega x} E(x) dx.$$

In Appendix A.2 it is shown that $\hat{e}(\omega) < \infty$ for $\omega > \lambda/\alpha$. Hence, for ω large enough and using that the q_n and $E(x)$ are non-negative to interchange the order of integration and summation (twice), we can rewrite (6.4) as:

$$\hat{e}(\omega) = \frac{1 - Q(\gamma^{-1})}{\alpha\omega^2} + \frac{\lambda}{\alpha\omega} \left(\hat{e}(\omega) - \sum_{n=1}^{\infty} q_n \hat{e}(\gamma^n \omega) \right).$$

This gives

$$\hat{e}(\omega) = \frac{1}{\alpha\omega - \lambda} \left(\frac{1 - Q(\gamma^{-1})}{\omega} - \lambda \sum_{n=1}^{\infty} q_n \hat{e}(\gamma^n \omega) \right). \quad (6.5)$$

Substituting this equation repeatedly into itself, applying partial fraction expansion at each step and using that $\hat{e}(\gamma^k \omega) \downarrow 0$ as $k \rightarrow \infty$ leads us to the following candidate solution:

$$\hat{e}(\omega) = \frac{1 - Q(\gamma^{-1})}{\omega} \sum_{i=0}^{\infty} \frac{e_i}{\gamma^i \alpha \omega - \lambda}, \quad (6.6)$$

where the e_i are constants to be determined. This representation will be justified by showing that it leads us to the (unique!) solutions to (6.5) and (6.4). Substituting (6.6) into (6.5) and equating the coefficients multiplying the terms $1/(\gamma^i \alpha \omega - \lambda)$ leads to:

$$\frac{e_i}{e_0} = \frac{1}{1 - \gamma^{-i}} \sum_{n=1}^i \gamma^{-n} q_n \frac{e_{i-n}}{e_0}, \quad i = 1, 2, 3, \dots, \quad (6.7)$$

$$e_0 = \left(1 + \sum_{n=1}^{\infty} \gamma^{-n} q_n \sum_{j=0}^{\infty} \frac{e_j/e_0}{\gamma^{j+n} - 1} \right)^{-1}. \quad (6.8)$$

We note that the ratios e_i/e_0 are non negative and can be computed recursively from (6.7). Then the normalizing constant $e_0 > 0$ can be computed from (6.8). From (6.7) it can be shown (by induction on i) that

$$e_i \leq \gamma^{-i} e_0, \quad i = 1, 2, \dots, \quad (6.9)$$

i.e., the e_i decay exponentially fast in i as $i \rightarrow \infty$. Therefore the right hand side of (6.6) certainly converges for $\omega > \lambda/\alpha$ and, from its construction, (6.6) is the solution to (6.5). By partial fraction (6.6) can be rewritten as:

$$\hat{e}(\omega) = \frac{1 - Q(\gamma^{-1})}{\lambda} \sum_{i=0}^{\infty} e_i \left(\frac{1}{\omega - \gamma^{-i} \lambda/\alpha} - \frac{1}{\omega} \right). \quad (6.10)$$

Inverting this LT gives:

$$E(x) = \frac{1 - Q(\gamma^{-1})}{\lambda} \sum_{i=0}^{\infty} e_i \left(e^{\gamma^{-i} (\lambda/\alpha)x} - 1 \right). \quad (6.11)$$

Using this in (6.1) we have

$$P_M = \frac{1/\lambda}{1/\lambda + E(M)} = \left(1 + (1 - Q(\gamma^{-1})) \sum_{i=0}^{\infty} e_i \left(e^{\gamma^{-i}(\lambda/\alpha)M} - 1 \right) \right)^{-1}. \quad (6.12)$$

Note that because of (6.9) and

$$\left(e^{\gamma^{-i}(\lambda/\alpha)M} - 1 \right) \sim \gamma^{-i}(\lambda/\alpha)M, \quad i \rightarrow \infty,$$

P_M can be computed efficiently from (6.12).

Remark 6.1 In particular cases we can find the coefficients e_i explicitly. For instance, when the reduction of the TCP window is always by a constant factor γ , i.e., $N \equiv 1$ (hence, $q_1 = 1$ and $Q(z) \equiv z$). Note that with $\gamma = 2$ we have TCP's most common window decrease factor (see Section 7 for more specific results in that case). From (6.7, 6.8) we get

$$\frac{e_i}{e_0} = \frac{\gamma^{-i}}{\prod_{j=1}^i (1 - \gamma^{-j})}, \quad i = 1, 2, \dots, \quad (6.13)$$

$$e_0 = \left(1 + \sum_{i=1}^{\infty} \gamma^{-i} \frac{e_i}{e_0} \right)^{-1}. \quad (6.14)$$

In this case we could have obtained these coefficients also in a direct way, without using (6.6), see Appendix A.3. There it is also shown that in this case

$$\sum_{i=0}^{\infty} e_i = \frac{1}{1 - \gamma^{-1}},$$

and, hence, from (6.12):

$$P_M = \left((1 - \gamma^{-1}) \sum_{i=0}^{\infty} e_i e^{\gamma^{-i}(\lambda/\alpha)M} \right)^{-1}. \quad (6.15)$$

7 Special case: single congestion signals and $\gamma = 2$

In this section we specify our results for the important particular case of TCP flow control with only one division of the window by a factor 2 at loss events. I.e., we take $\gamma = 2$ and $N \equiv 1$ ($q_1 = 1$, and $q_n = 0, n = 2, 3, \dots$) in the model with peak rate limitation, see [14] for a similar model. In Section 9.2 we compare the results from this particular case of our model to measurements from the Internet. We worked with long distance connections where congestion signals rarely appear in batches. From (4.4) and (4.5) we obtain the expressions for the first two moments of the window size distribution:

$$E[W] = \frac{2\alpha}{\lambda}(1 - P_M).$$

$$E[W^2] = \frac{8\alpha[2\alpha(1 - P_M) - \lambda P_M M]}{3\lambda^2},$$

where P_M is given by (6.15) with $\gamma = 2$. The throughput of TCP can be obtained from Equation (4.1). The distribution function itself or the complementary distribution function $\overline{F}(x)$ is computed successively on the intervals $[M/2^k, M/2^{k-1}]$, $k = 1, 2, \dots$ using (5.5) with $\gamma = 2$. Recursion (5.6) reduces to

$$c_{i+1}^{(k)} = \frac{c_i^{(k-1)}}{1 - 2^i}, \quad i = 1, \dots, k - 1,$$

and $c_1^{(k)}$ is given by (5.7). When $M = \infty$, the distribution function is given by (5.12) with

$$c_i = \frac{1}{1 - 2^i} c_{i-1}, \quad i = 1, 2, \dots,$$

and c_0 is given by (5.11).

8 General congestion limitation model

In the case of congestion limitation it seems unrealistic to assume that the number of congestion signals in the batch that is generated when reaching the maximum transmission rate has the same distribution as the size of the batches generated by the Poisson process. Let us therefore assume that the number of congestion signals that result from reaching the maximum transmission rate is distributed as the non negative discrete random variable $N^{(M)}$ having p.g.f.

$$Q^{(M)}(z) = \sum_{n=1}^{\infty} z^n q_n^{(M)} \quad |z| \leq 1.$$

Instead of (2.5) we then have for $0 < x < M$:

$$-\alpha \frac{d}{dx} \overline{F}^{\text{cl}}(x) = \lambda \left(\overline{F}^{\text{cl}}(x) - \sum_{n=1}^{\infty} q_n \overline{F}^{\text{cl}}(\min(\gamma^n x, M)) \right) + \lambda b^{(M)} H(x), \quad (8.1)$$

where

$$H(x) := \mathbf{P} \left\{ N^{(M)} \geq \frac{\ln(M) - \ln(x)}{\ln(\gamma)} \right\},$$

$$b^{(M)} := -\frac{\alpha}{\lambda} \frac{d}{dy} \overline{F}^{\text{cl}}(y) \Big|_{y=M-}.$$

Note that $H(x)$ is a non negative, non decreasing step function of the variable x , constant on the intervals $\gamma^{-k} M < x \leq \gamma^{-k+1} M$, $k = 1, 2, \dots$, with $H(M) = 1$:

$$H(x) = h_k := \mathbf{P} \left\{ N^{(M)} \geq k \right\}, \quad \gamma^{-k} M \leq x < \gamma^{-k+1} M.$$

8.1 The moments

Similar to Section 4 we find the following recursion on the moments after taking Laplace Transforms in (8.1):

$$\mathbf{E} \left[\left(W^{\text{cl}} \right)^k \right] = \frac{k\alpha \mathbf{E} \left[\left(W^{\text{cl}} \right)^{k-1} \right] - b^{(M)} \lambda M^k (1 - Q^{(M)}(\gamma^{-k}))}{\lambda (1 - Q(\gamma^{-k}))}.$$

Note that if $Q^{(M)}(z) = Q(z)$ this recursion indeed reduces to (2.8).

8.2 The distribution function

Defining, for $k = 1, 2, \dots$,

$$\overline{F}_k^{\text{cl}}(x) := \overline{F}^{\text{cl}}(x), \quad \frac{M}{\gamma^k} \leq x < \frac{M}{\gamma^{k-1}},$$

we find by the same techniques as in Section 5:

$$\begin{aligned} \overline{F}_1^{\text{cl}}(x) &= b^{(M)} \left(e^{\frac{\lambda}{\alpha}(M-x)} - 1 \right), \\ \overline{F}_k^{\text{cl}}(x) &= -b^{(M)} h_k + \left(\overline{F}_{k-1}^{\text{cl}}(M/\gamma^{k-1}) + b^{(M)} h_k \right) e^{\frac{\lambda}{\alpha} \left(\frac{M}{\gamma^{k-1}} - x \right)} \\ &\quad - \frac{\lambda}{\alpha} e^{-\frac{\lambda}{\alpha} x} \int_{u=x}^{M/\gamma^{k-1}} e^{\frac{\lambda}{\alpha} u} \sum_{n=1}^{k-1} q_n \overline{F}_{k-n}^{\text{cl}}(\gamma^n u) du. \end{aligned}$$

This leads to:

$$\overline{F}_k^{\text{cl}}(x) = b^{(M)} \left(-d_0^{(k)} + \sum_{i=1}^{k-1} d_i^{(k)} e^{-\frac{\lambda}{\alpha} \gamma^{i-1} x} \right).$$

The coefficients $d_i^{(k)}$ are given by:

$$\begin{aligned} d_0^{(k)} &= h_k + \sum_{n=1}^{k-1} q_n d_0^{(k-n)}, \quad k > 1, \\ d_i^{(k)} &= \frac{1}{1 - \gamma^{i-1}} \sum_{n=1}^{i-1} q_n d_{i-n}^{(k-n)}, \quad k > i > 1, \end{aligned}$$

and for $k > 1$:

$$\begin{aligned} d_1^{(k)} &= \left(-d_0^{(k-1)} + \sum_{i=1}^{k-2} d_i^{(k-1)} e^{\frac{\lambda}{\alpha} M/\gamma^{k-i}} + h_k \right) e^{\frac{\lambda}{\alpha} M/\gamma^{k-1}} \\ &\quad + \sum_{n=1}^{k-1} q_n e^{\frac{\lambda}{\alpha} M/\gamma^{k-n-1}} \left(d_0^{(k-n)} - \sum_{i=1}^{k-n-1} d_i^{(k-n)} \frac{1}{1 - \gamma^{i-1+n}} e^{-\frac{\lambda}{\alpha} M/\gamma^{k-i-2n}} \right). \end{aligned}$$

Note that the $d_0^{(k)}$ are all non negative, but that the signs of the $d_i^{(k)}$ for $i > 0$ alternate.

8.3 The constant $b^{(M)}$

Note that, similar to (5.8) we find from $\overline{F}^{\text{cl}}(0) = 1$:

$$\frac{1}{b^{(M)}} = \lim_{k \rightarrow \infty} \left(-d_0^{(k)} + \sum_{i=1}^{k-1} d_i^{(k)} e^{-\frac{\lambda}{\alpha} M / \gamma^{k-i}} \right). \quad (8.2)$$

However, for computational purposes, we again prefer to translate the model into a peak rate limitation type of model. Therefore, consider a peak rate limitation model in which congestion signal batches arrive according to a Poisson process with rate λ but, different from the model in Section 2.1, with the distribution of the batch size depending on whether the transmission rate is below or at the maximum level M . The batch size has p.g.f. $Q(z)$ if the rate is below M , otherwise it has p.g.f. $Q^{(M)}(z)$. Similar to (2.9) we have:

$$b^{(M)} = \frac{P_M^{(M)}}{1 - P_M^{(M)}},$$

where $P_M^{(M)}$ is the probability of being at the maximum transmission rate M . For $0 < x < M$ (as we shall see it is convenient not to include $x = M$), let the functions $E(x)$ and $E_n(x)$, $n = 1, 2, \dots$, be defined as in Section 6. Note that as long as the process is below the maximum level M , it behaves exactly as the ordinary peak rate limitation model of Section 2.1. Therefore, for $0 < x < M$ the functions $E(x)$ and $E_n(x)$ are exactly as we found in Section 6, see (6.11) and (6.3). This is not true for $x = M$ and to avoid confusion we write $E^{(M)}(M)$ instead of $E(M)$ for the return time to level M in the present model. Of course,

$$P_M^{(M)} = \frac{1}{1 + \lambda E^{(M)}(M)}.$$

Similar to (6.2) we have:

$$E^{(M)}(M) = \sum_{n=1}^{\infty} q_n^{(M)} E_n(M).$$

And using (6.3) and (6.11) we find:

$$\begin{aligned} E^{(M)}(M) &= \frac{(1 - Q^{(M)}(\gamma^{-1})) M}{\alpha} \left[1 - (1 - Q(\gamma^{-1})) \sum_{i=0}^{\infty} e_i \right] \\ &\quad + \frac{1 - Q(\gamma^{-1})}{\lambda} \sum_{i=0}^{\infty} \gamma^i e_i \sum_{n=1}^{\infty} q_n^{(M)} \left(e^{\gamma^{-i}(\lambda/\alpha)M} - e^{\gamma^{-i-n}(\lambda/\alpha)M} \right). \end{aligned}$$

9 Model validation

In this section we compare measurements from long distance and long life TCP connections with the results of Section 7 ($N \equiv 1, \gamma = 2$, peak rate limitation). Comparison of real measurements with the model with clustered (batch) arrivals of congestion signals is a topic of current research, see also Section 10.

Due to the large number of hops and the multiplexing of exogenous traffic in network routers, the Poisson loss process assumption is expected to hold on long distance connections [14]. Our TCP receivers implement the delay ACK mechanisms and our TCP senders increase their window in the congestion avoidance mode by approximately one packet every window's worth of ACKs. Thus, we take α equal to $1/(2RTT)$ [18]. First, we show theoretically how the window size is distributed in the stationary regime. Second, we compare our results to measurements from the Internet.

9.1 Numerical results

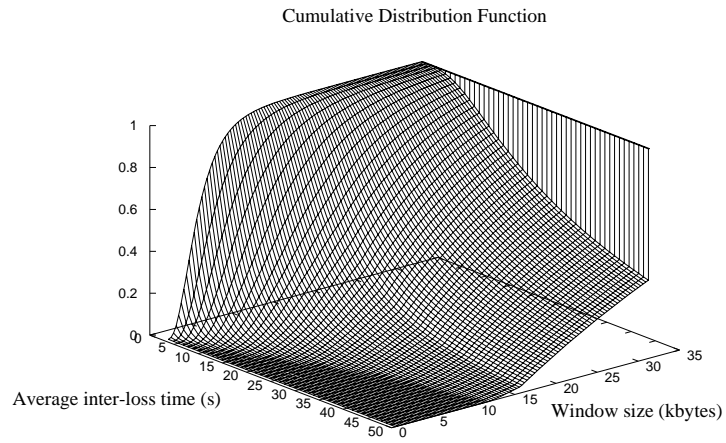
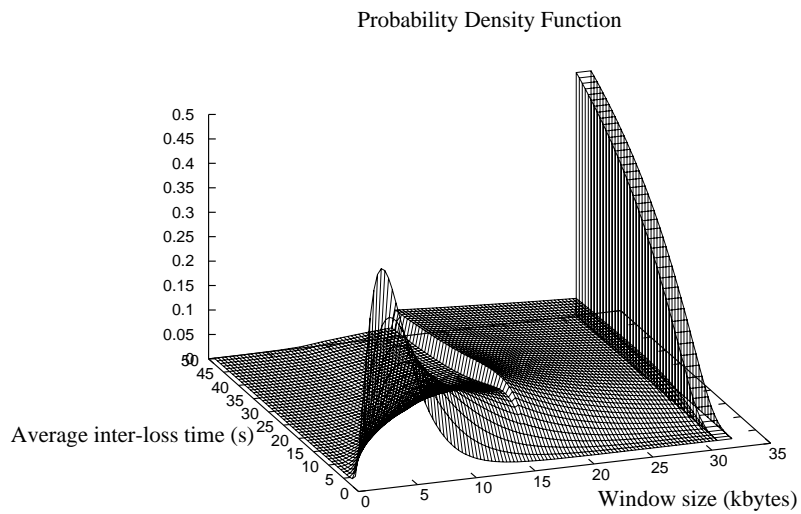
Consider the case of a long TCP connection with packets of size 1460 bytes and a constant RTT of one second. Using the results of Section 7, we plot the cumulative distribution function $F(x)$ of the window size and its probability density function $f(x)$ for a range of values for the loss intensity λ (or rather, for the mean inter-loss time $s = 1/\lambda$). Two values of M are considered. First, we let the congestion window be limited by a receiver window of 32 kbytes. Then we consider the case where the window is not limited and therefore continues to grow linearly until a loss occurs. The numerical results are presented in Figures 3, 4, 5 and 6. For the case $M = 32$ kbytes, we computed the distribution function successively for the intervals $[M/2, M]$, $[M/4, M/2]$ and so on. When computing P_M we truncate the infinite series in (6.15). In the case of an unlimited window, we also truncate the infinite series in (5.12). As discussed previously, these infinite series converge fast. We choose the number of terms of these series large enough to get a negligible error.

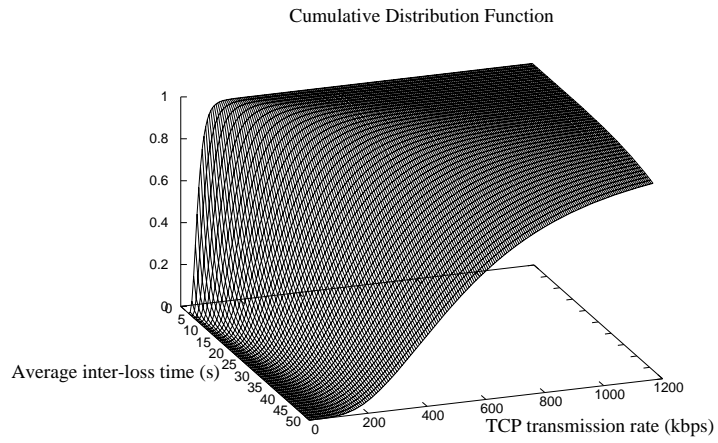
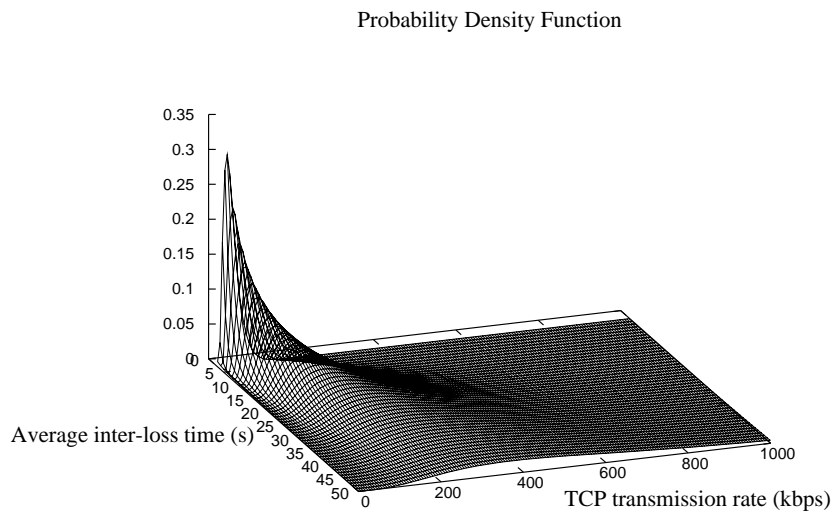
When $M = 32$ kbytes, the discontinuities of $F(x)$ and $f(x)$ at $x = M$ and $x = M/2$, respectively, are clearly illustrated in Figures 3 and 4. The discontinuities are most noticeable for large inter-loss times. The discontinuity of $F(x)$ is also depicted in Figure 4 by plotting a pulse for $f(x)$ at $x = M$ such that its area is equal to P_M . When $M = \infty$, the density function exhibits neither pulses nor discontinuities (Figure 6).

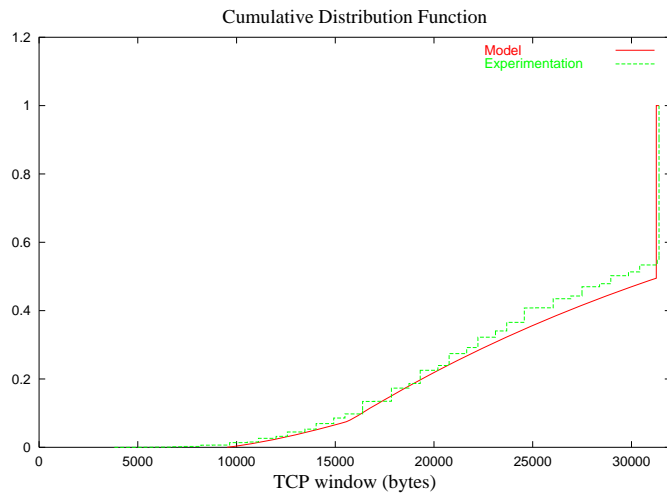
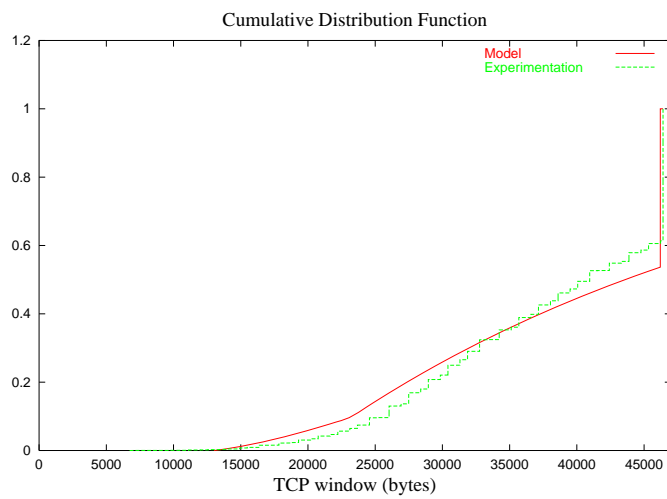
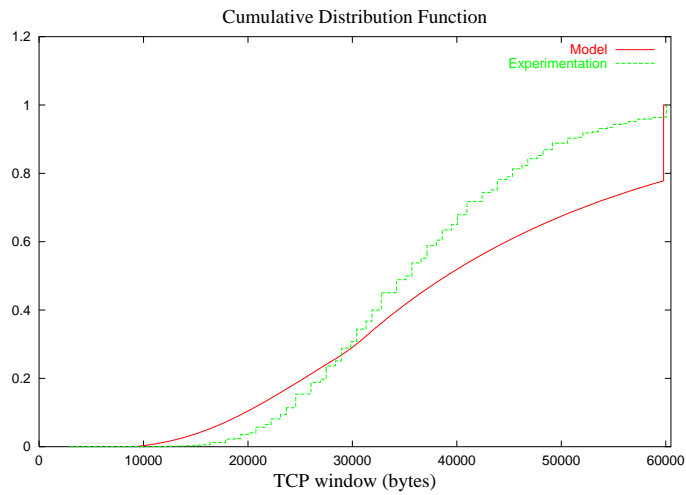
9.2 Experimental results

Our experimental testbed consists of a long life and long distance TCP connection between INRIA Sophia Antipolis (France) and Michigan State University (US). The TCP connection is fed at INRIA by an infinite amount of data. The New Reno version of TCP [7] is used for data transfer. We change the socket buffer at the receiver in order to account for different values of M . We considered three values of M : 32, 48 and 64 kbytes. For every value of M , we ran the TCP connection for approximately one hour and we registered the trace of the connection using the `tcpdump` tool of LBNL [12]. We developed a tool that analyzes the trace of the connection and that detects the times at which the window is reduced. This tool gives also the average RTT of the connection and the statistics of the window per RTT. We plotted for the three values of M , the distribution of the window size from measurements and that given by our model. The results are plotted in Figures 7, 8 and 9.

When M is small, we observe a good match between the measured distribution and the one resulting from our model. For larger values of M , the difference between the two increases. In particular, as M increases, the measured probability density concentrates around the average window size. This deviation can be explained from the measured inter-loss time distribution.

Figure 3: Limited receiver window: $M = 32$ kbytesFigure 4: Limited receiver window: $M = 32$ kbytes

Figure 5: Unlimited receiver window: $M = \infty$ Figure 6: Unlimited receiver window: $M = \infty$

Figure 7: Receiver window $M = 32$ kbytesFigure 8: Receiver window $M = 48$ kbytesFigure 9: Receiver window $M = 64$ kbytes

In Figure 10, we plot this distribution for $M = 32$. This distribution is in agreement with an exponential law, resulting in a good match between the model and the measurements. Figures 11 and 12 show the measured distributions for the other two values of M . We observe that the loss process is no longer Poisson, but closer to a deterministic process. Small inter-loss times are less frequent as M increases and the tail of the distribution also becomes less important (although it still looks like the tail of an exponential distribution). This results in a degradation of the correspondence between our model and the measurements.

One explanation of the deviation of the loss process from a Poisson process for larger values of M is the following. A true Poisson loss process implies that the time until the next loss event is independent of the past. This is the case when the congestion of the network is dominated by the exogenous traffic and not dependent on the measured connection. I.e., it is the case when the measured connection's share of the available bandwidth on the path is small compared to that of the exogenous traffic. A small M limits the bandwidth share of our connection and limits its impact on the network, resulting in a loss process close to Poisson. However for large M , the measured connection achieves a larger share and thus contributes more to the congestion of network routers. When it reduces its window, the state of the network changes and becomes under-loaded. For a certain amount of time the occurrence of a new congestion is less likely. When the network is again more heavily loaded the next congestion signal is likely to appear soon. This explains why we observe a low density for small inter-loss times (Figures 12), then a peak in the middle followed by an exponential decay.

In summary, our model leads to accurate results when the times between losses are exponentially distributed. However, in situations where the congestion in the network is largely due to the TCP connection under consideration, the loss process is close to a deterministic process and a simple heuristic as that proposed in [3, 18] can be used to approximate the achieved throughput.

10 Conclusions and future research

We studied additive-increase multiplicative-decrease flow control mechanisms under the assumption that congestion signals arrive in batches according to a Poisson process. As highlighted in [14], the model can be reformulated as an M/G/1 queuing problem with service time dependent on system workload. We tried to keep the model as general as possible in order to account for a wide range of congestion control strategies. We derived closed form expressions for the moments as well as the distribution of the transmission rate. For the case of single congestion signals, we compared our results to measurements from TCP connections over the Internet. From our experiments, we concluded that our model with single congestion signals leads to accurate results when the times between losses are close to being exponentially distributed.

Currently, we are working on the validation of our model with clustered congestion signals. Internet measurements have shown that on some paths (especially short distance ones) the loss process exhibits a high degree of burstiness. We also study the extension of the analysis

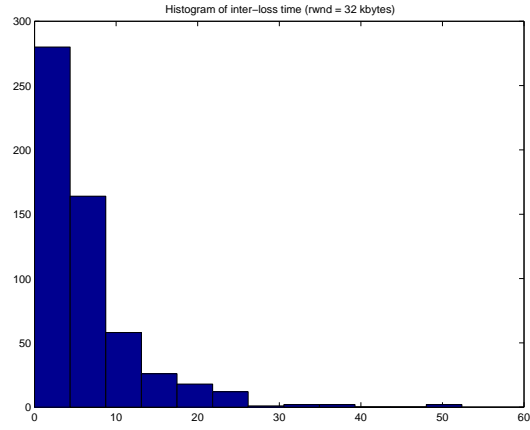


Figure 10: Case of $M = 32$ kbytes

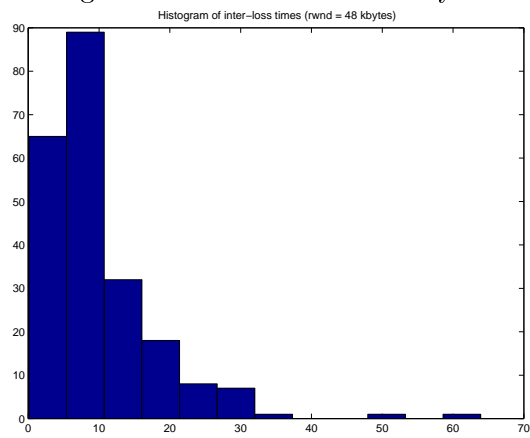


Figure 11: Case of $M = 48$ kbytes

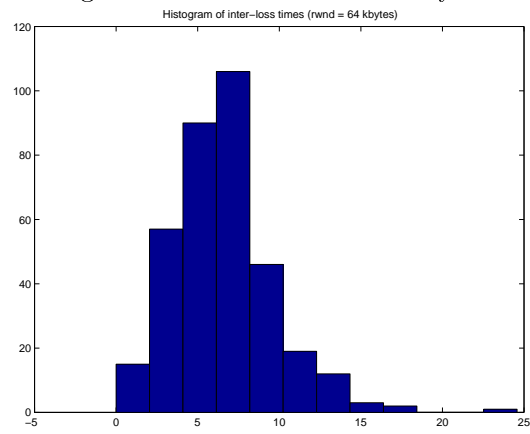


Figure 12: Case of $M = 64$ kbytes

to more general inter-loss processes, in particular to Markov Modulated Poisson Processes.

Acknowledgement

The authors would like to thank V.M. Abramov for comments that improved the presentation of the paper, as well as the department of Mathematics of Michigan State University and in particular L.B. Freidovich for providing us with a computer account for our Internet experiments.

Appendix

A.1 Uniqueness of $E(x)$

If $\tilde{E}(x)$ is a second solution to (6.4) then $D(x) := E(x) - \tilde{E}(x)$ satisfies:

$$D(x) = \frac{\lambda}{\alpha} \sum_{n=1}^{\infty} q_n \int_{y=\gamma^{-n}x}^x D(y) dy, \quad x \geq 0,$$

hence,

$$\begin{aligned} |D(x)| &\leq \frac{\lambda}{\alpha} \sum_{n=1}^{\infty} q_n \int_{y=\gamma^{-n}x}^x |D(y)| dy \\ &\leq \frac{\lambda}{\alpha} \int_{y=0}^x |D(y)| dy. \end{aligned} \tag{A.1}$$

Define the function $h(x)$, $x \geq 0$, by

$$\begin{aligned} \int_{y=0}^x |D(y)| dy &=: e^{\frac{\lambda}{\alpha}x} h(x), \\ |D(x)| &= \frac{\lambda}{\alpha} e^{\frac{\lambda}{\alpha}x} h(x) + e^{\frac{\lambda}{\alpha}x} \frac{d}{dx} h(x). \end{aligned}$$

Substitution into (A.1) gives

$$\frac{d}{dx} h(x) \leq 0,$$

Obviously, from its definition above, $h(x) \geq 0$ and $h(0) = 0$, hence, $h(x) = 0$ for all $x \geq 0$. This proves that $D(x) \equiv 0$.

A.2 Existence of $\hat{e}(\omega)$

Using that $E(y)$ is non negative for $y \geq 0$ it follows from (6.4) that

$$E(x) \leq \frac{(1 - Q(\gamma^{-1}))x}{\alpha} + \frac{\lambda}{\alpha} \int_{y=0}^x E(y) dy.$$

Writing

$$\int_{y=0}^x E(y) dy =: e^{\frac{\lambda}{\alpha}x} h(x),$$

gives

$$\frac{d}{dx}h(x) \leq \frac{(1 - Q(\gamma^{-1}))x}{\alpha} e^{-\frac{\lambda}{\alpha}x},$$

hence,

$$0 \leq h(x) \leq \frac{1 - Q(\gamma^{-1})}{\alpha} \left[\left(\frac{\alpha}{\lambda} \right)^2 \left(1 - e^{-\frac{\lambda}{\alpha}x} \right) - \frac{\alpha}{\lambda} x e^{-\frac{\lambda}{\alpha}x} \right].$$

Therefore $\hat{e}(\omega) < \infty$ for $\omega > \lambda/\alpha$.

A.3 Direct derivation of P_M for single congestion signals

Consider the case where $N \equiv 1$. We show a direct approach to find the coefficients e_i without using the ‘candidate’ (6.6). We shall need the following identities (which hold for $\gamma > 1$):

$$\left(1 + \sum_{i=1}^{\infty} \frac{\gamma^{-2i}}{\prod_{j=1}^i (1 - \gamma^{-j})} \right)^{-1} = \sum_{k=0}^{\infty} \frac{\gamma^{-k}}{\prod_{l=1}^k (1 - \gamma^l)}, \quad (\text{A.2})$$

$$\sum_{i=0}^k \frac{1}{\prod_{m=1}^{k-i} (1 - \gamma^m) \prod_{n=1}^i (1 - \gamma^{-n})} = 1. \quad (\text{A.3})$$

Equation (6.5) becomes:

$$\hat{e}(\omega) = \frac{1}{\alpha\omega - \lambda} \left(\frac{1 - \gamma^{-1}}{\omega} - \lambda \hat{e}(\gamma\omega) \right),$$

and substituting this equation repeatedly into itself we find:

$$\hat{e}(\omega) = \frac{1 - \gamma^{-1}}{\omega} \sum_{k=0}^{\infty} \frac{(-\lambda\gamma^{-1})^k}{\prod_{j=0}^k (\gamma^j \alpha\omega - \lambda)}. \quad (\text{A.4})$$

Note that the infinite sum converges absolutely for all $\omega \neq 0$ and $\omega \neq \gamma^{-j}\lambda/\alpha$, for $j = 0, 1, 2, \dots$, because the k -th term is of the order of $\gamma^{-k^2/2}$ when $k \rightarrow \infty$. By partial fraction expansion we have:

$$\begin{aligned} \hat{e}(\omega) &= (1 - \gamma^{-1}) \sum_{k=0}^{\infty} \gamma^{-k} \left[-\frac{1}{\lambda\omega} + \frac{\alpha}{\lambda} \sum_{i=0}^k \frac{\left(\prod_{j=0, j \neq i}^k (1 - \gamma^{j-i}) \right)^{-1}}{\alpha\omega - \gamma^{-i}\lambda} \right] \\ &= -\frac{1}{\lambda\omega} + (1 - \gamma^{-1}) \frac{\alpha}{\lambda} \sum_{k=0}^{\infty} \gamma^{-k} \sum_{i=0}^k \frac{\left(\prod_{j=0, j \neq i}^k (1 - \gamma^{j-i}) \right)^{-1}}{\alpha\omega - \gamma^{-i}\lambda}. \end{aligned} \quad (\text{A.5})$$

(By convention the empty product equals 1.) Using the absolute convergence of the infinite series we have:

$$E(x) = -\frac{1}{\lambda} + \frac{1 - \gamma^{-1}}{\lambda} \sum_{k=0}^{\infty} \gamma^{-k} \sum_{i=0}^k \left(\prod_{j=0, j \neq i}^k (1 - \gamma^{j-i}) \right)^{-1} \exp\left(\frac{\lambda x}{\alpha\gamma^i}\right). \quad (\text{A.6})$$

Using (A.2) we obtain the coefficients e_i in (6.13) and (6.14). If we use (A.3) we can also show that in this case

$$\sum_{i=0}^{\infty} e_i = \frac{1}{1 - \gamma^{-1}},$$

and, hence, from (6.12) we obtain (6.15).

References

1. THE ATM FORUM TECHNICAL COMMITTEE. Traffic Management Specification. Version 4.0, 95-0013R8, 1995.
2. V.M. ABRAMOV. Investigation of a queueing system with service depending on queue length (in Russian). Donish (pub.), Dushanbe (1991).
3. E. ALTMAN, K. AVRACHENKOV, C. BARAKAT. TCP in presence of bursty losses. *ACM SIGMETRICS* (2000).
4. E. ALTMAN, K. AVRACHENKOV, C. BARAKAT. A stochastic model of TCP/IP with stationary random losses. *ACM SIGCOMM* (2000). Also available as INRIA Research Report RR-3824.
5. D. CHOI, C. KNESSL, C. TIER. A queueing system with queue length dependent service times, with applications to cell discarding in ATM networks. *J. Appl. Math. Stoch. Anal.* 12 (1999), no. 1, 35–62.
6. S. FLOYD. TCP and Explicit Congestion Notification. *ACM Computer Communication Review* (1994).
7. S. FLOYD, T. HENDERSON. The NewReno Modification to TCP's Fast Recovery Algorithm. *RFC 2582* (1999).
8. J. HEINANEN, T. FINLAND, F. BAKER, W. WEISS, J. WROCLAWSKI. Assured Forwarding PHB Group. *RFC 2597* (1999).
9. V. JACOBSON. Congestion avoidance and control. *ACM SIGCOMM* (1988).
10. C. KNESSL, B.J. MATKOWSKY, Z. SCHUSS, C. TIER. Asymptotic Analysis of a state-dependent M/G/1 queueing system. *SIAM J. Appl. Math.* 46, no. 3 (1986), 483–505.
11. C. KNESSL, C. TIER, B.J. MATKOWSKY, Z. SCHUSS. A state-dependent GI/G/1 queue. *Eur. J. Appl. Math.* 5, no. 2 (1994), 217–241.
12. LBNL's tcpdump tool, available at <http://www-nrg.ee.lbl.gov>.

13. M. MATHIS, J. SEMKE, J. MAHDAVI, T. OTT. The Macroscopic Behavior of the TCP Congestion Avoidance Algorithm. *ACM Computer Communication Review*, June 1997.
14. V. MISRA, W.-B. GONG, D. TOWSLEY. Stochastic differential equation modeling and analysis of TCP-window size behavior. *Performance'99, Istanbul (Turkey)*. Oct. 1999 (available at <http://www-net.cs.umass.edu/papers/papers.html>).
15. A. MISRA, T. OTT, J. BARAS. The Window Distribution of Multiple TCPs with Random Queues. *IEEE GLOBECOM*, Dec 1999.
16. S.C. NIU. On queues with dependent interarrival and service times. *Naval Res. Logist. Quart.* 28, no. 3 (1981), 497–501.
17. T.J. OTT, J.H.B. KEMPERMAN, M. MATHIS. The stationary behavior of ideal TCP congestion avoidance. available at <ftp://ftp.telecordia.com/pub/tjo/TCPwindow.ps> (August 1996).
18. J. PADHYE, V. FIROIU, D. TOWSLEY, J. KUROSE. Modeling TCP throughput: a simple model and its empirical validation. *ACM SIGCOMM*, Sep. 1998.
19. M. POSNER. Single-server queues with service time dependent on waiting time. *Operations Res.* 21 (1973), 610–616.
20. S. SAHU, D. TOWSLEY, J. KUROSE. A Quantitative Study of Differentiated Services for the Internet. *IEEE GLOBECOM*, Dec. 1999.
21. W. STEVENS. TCP Slow Start, Congestion Avoidance, Fast Retransmit, and Fast Recovery Algorithms. *RFC 2001* (Jan. 1997).
22. H.C. TIJMS. *Stochastic Models — An Algorithmic Approach*. Wiley, Chichester (1994).