

## **BUFFER OVERFLOW ASYMPTOTICS FOR A BUFFER HANDLING MANY TRAFFIC SOURCES**

COSTAS COURCOUBETIS,\* *University of Crete*  
RICHARD WEBER,\*\* *University of Cambridge*

### **Abstract**

As a model for an ATM switch we consider the overflow frequency of a queue that is served at a constant rate and in which the arrival process is the superposition of  $N$  traffic streams. We consider an asymptotic as  $N \rightarrow \infty$  in which the service rate  $Nc$  and buffer size  $Nb$  also increase linearly in  $N$ . In this regime, the frequency of buffer overflow is approximately  $\exp(-NI(c, b))$ , where  $I(c, b)$  is given by the solution to an optimization problem posed in terms of time-dependent logarithmic moment generating functions. Experimental results for Gaussian and Markov modulated fluid source models show that this asymptotic provides a better estimate of the frequency of buffer overflow than ones based on large buffer asymptotics.

ATM SWITCHES; BUFFER OVERFLOW ASYMPTOTICS; EFFECTIVE BANDWIDTHS; LARGE DEVIATIONS; MARKOV MODULATED FLUID

AMS 1991 SUBJECT CLASSIFICATION: PRIMARY 60K30  
SECONDARY 60F10; 60K25; 68M20; 90B10; 90B22

### **1. Switches handling many bursty sources**

In a high speed data communications network, such as one operating according to the ideas of ATM (asynchronous transfer mode), data is packaged into cells of fixed size which are transmitted between switches along high capacity links. The traffic sources are of various types, such as voice, video and file transfer, and they are bursty, in the sense that the rate at which cells are produced by a source is not constant, but fluctuates around a mean rate. Switches are buffered, as a safeguard against those occasions when the total rate at which cells enter the buffer from all sources that are routed through that switch exceeds the rate at which they can be served by the output links. Occasionally, the buffer will not be large enough and cell loss will occur. This should happen rarely. Since the number of sources that are routed through a switch is large, there is a statistical multiplexing which reduces the probability of buffer overflow. The idea is that when some sources are producing cells at above their average rates other sources will be producing cells at below their average rates.

---

Received 9 December 1994; revision received 31 May 1995.

\* Postal address: Department of Computer Science, University of Crete, P.O. Box 1470 Heraklion, 71110 Greece.

\*\* Postal address: Statistical Laboratory, 16 Mill Lane, Cambridge CB2 1SB, UK.

A measure of Quality of Service (QoS) is the cell loss rate; this rate should be very small, typically of the order of  $10^{-6}$  to  $10^{-10}$ . The cell loss rate is difficult to measure, but is important to estimate, both to decide questions of call acceptance and to identify paths in the network that are heavily or lightly loaded for the purpose of call routing. Except for very simple models it is impossible for a queueing theory analysis to evaluate the savings in bandwidth due to statistical multiplexing or to estimate the cell loss rate. For this reason researchers have focused their attention on analysis based on asymptotics and upon on-line measurements of the cell loss rate. One such asymptotic has been developed for systems with large buffers. In this paper we study an asymptotic for a large numbers of sources, and show that it can be used more successfully to estimate the cell loss rate.

1.1. *The model.* The behavior of a switch in an ATM network can be modelled as a queue with a constant service rate and a buffer for  $B$  cells. We shall suppose that time is discretised into epochs,  $n=1, 2, \dots$ , and that during epoch  $n$  the cell service rate and arrival rate are both constant, and equal to  $C$  and  $X_n$  cells per epoch respectively. The workload at the start of epoch  $n$  is denoted  $W_n$  and

$$W_{n+1} = \max\{0, \min[(W_n + X_n - C), B]\}.$$

The minimization in the above reflects the fact that cells are lost when the buffer is full, and the maximization ensures that the queue cannot become negative. We shall assume that  $\{X_n\}$  is an ergodic process, and that  $E[X_n] < C$ . This condition implies that  $\{W_n\}$  is ergodic.

Suppose that  $\{X_n\}$  is the superposition of sources of  $M$  different types, with there being  $N\rho_i$  sources of type  $i$ ,  $i=1, \dots, M$ . Thus the traffic is defined by the parameter  $\rho = (\rho_1, \dots, \rho_M)$  and the scaling parameter  $N$ . Let  $b$  and  $c$  denote respectively the amounts of buffer space and bandwidth per source, so that  $B = Nb$  and  $C = Nc$ .

The cell loss rate can be expressed as  $L(c, b, N) = E[(W_n + X_n - C - B, 0)^+]$ , when  $W_n$  and  $X_n$  have their stationary values. A related measure is the proportion of time the buffer is full, which we denote  $\Phi(c, b, N) = P(W_n = B)$ . Alternatively, we might suppose the buffer is infinite and study the proportion of time that the content is above level  $B$ , say  $Q(c, b, N)$ . Note that  $L(c, b, N)$ ,  $\Phi(c, b, N)$  and  $Q(c, b, N)$  are all functions of  $\rho$ , but to lighten the notation this dependence is not shown.

1.2. *The large  $N$  asymptotic.* The large  $N$  asymptotic with which we are concerned takes the following form:

$$(1) \quad \Phi(c, b, N) = \exp(-NI(c, b) + g_1(c, b, N)),$$

where  $\lim_{N \rightarrow \infty} g_1(c, b, N)/N = 0$ . Both  $I(c, b)$  and  $g_1(c, b, N)$  depend on  $\rho$  but as we have done above this is suppressed to lighten the notation. The rate  $I(c, b)$  is found as the solution to an optimization problem posed in terms of time dependent logarithmic moment generating functions. Asymptotic (1) can be compared with a well-known asymptotic in  $B$ ,

$$(2) \quad \Phi(c, b, N) = \exp(-NbH(c) + g_2(c, b, N)),$$

where  $\lim_{b \rightarrow \infty} g_2(c, b, N)/b = 0$ . This is an asymptotic that has been studied in a long series of papers, by Courcoubetis and Weber (1995), de Veciana and Walrand (1994), Elwalid and Mitra (1993), Kelly (1991), Kesidis *et al.* (1993) and Whitt (1993). Either (1) or (2) may be used to estimate  $\Phi$ , as  $\exp(-NI(c, b))$  or  $\exp(-NbH(c))$  respectively. In Section 2 we show that  $I(c, b)/b \rightarrow H(c)$  as  $b \rightarrow \infty$  and so if the buffer space per source is large then either of these estimates might be used.

However, we find that it is usually better to base an estimate of  $\Phi$  on (1) rather than (2). For one thing, to base an estimate upon the asymptotic for large  $N$  takes more account of what really happens in practice. In the output link of an actual ATM switch, the number of sources is of the order  $10^3$ – $10^4$  and so it is correct to apply a result that assumes  $N$  is large. One cannot say the same for the approach that assumes the buffer is large; in many ATM switch designs, the total buffer space may be moderate, say 100–200 cells. This means that in actuality there is a much smaller amount of buffering per source than is supposed by use of the estimate based on the large  $b$  asymptotic.

The utilization of the switch by the bursty traffic also has an impact on the relative merits of the approaches. In an ATM network, some of the real-time traffic, such as voice and video, will be assigned high priority and some, such as file transfer and interactive traffic, can be delayed and subject to flow control. The real-time traffic will only utilize a proportion of the bandwidth. If we are concerned to model the high-priority traffic  $\{X_n\}$  can be considered to be a model for this traffic alone. Simulations show that when the utilization of the switch is low, say  $X_n = (2/3)C$ , the estimation of the overflow rate using the large  $N$  asymptotic remains accurate, whereas the estimate based on the large  $b$  asymptotic is very poor, overestimating the overflow rate by as much as  $10^5$ .

Section 2 begins with a proof of the large  $N$  asymptotic. We discuss the form of  $I(c, b)$  in cases when  $b$  is small or large, and the implications of expressing a quality of service constraint as the requirement that  $I(c, b)$  should exceed some target value.

In Section 3 we calculate  $I(c, b)$  for a Markov modulated fluid model of a source and observe from some calculations with typical values that the large  $N$  asymptotic does a much better job of estimating the cell loss rate for this model than does the large  $b$  asymptotic. Our approach reproduces more easily some of the results of Weiss that were obtained by a more refined analysis.

In Section 4 we consider a model of a source as an autoregressive Gaussian process. Again, the large  $N$  asymptotic does a better job of estimating the cell loss rate than does the large  $b$  asymptotic. In the special case where the parameters of the autoregressive process are chosen so that the source model becomes relatively unbursty then the large  $b$  asymptotic underestimates the cell loss rate, whereas the large  $N$  asymptotic continues to overestimate it and therefore leads to conservative decisions. The analysis in this section confirms and explains observations of Choudhury *et al.* (1994a), (1994b).

Section 5 concludes with a discussion of on-line estimation and of asymptotics for traffic whose composition is unknown.

**2. The large  $N$  asymptotic**

Both the large  $N$  and large  $b$  asymptotics are expressed in terms of logarithmic moment generating functions. Suppose  $\{Y_1^i, Y_2^i, \dots\}$  denotes cells generated in epochs 1, 2, ..., by a traffic source of type  $i$ , and define

$$\varphi_m^i(\theta) = \frac{1}{m} \log E \left[ \exp \left( \theta \sum_{n=1}^m Y_n^i \right) \right].$$

It is well-known that  $\varphi_m^i(\theta)$  is a convex increasing function of  $\theta$ . We shall also suppose that the asymptotic logarithmic moment generating function exists, defined as

$$\varphi^i(\theta) = \lim_{m \rightarrow \infty} \varphi_m^i(\theta).$$

We write  $\varphi_m(\theta) = \sum_i \rho_i \varphi_m^i(\theta)$  and  $\varphi(\theta) = \lim_{m \rightarrow \infty} \varphi_m(\theta)$ .

The large  $b$  asymptotic of (2) has its rate given by

$$(3) \quad H(c) = \sup[\delta : c \geq \varphi(\delta)/\delta].$$

Note that  $H(c)$  does not depend on  $N$ , but only on proportions of the different types of traffic. Since the sources are independent, and  $\varphi^i(\delta)/\delta$  is increasing in  $\delta$ ,

$$(4) \quad H(c) \geq \delta \Leftrightarrow c \geq \sum_i \rho_i \frac{\varphi^i(\delta)}{\delta},$$

and the quantity  $\varphi^i(\delta)/\delta$  is identified as the effective bandwidth of source type  $i$ .

Exactly the same asymptotic holds for  $L(c, b, N)$ ,  $\Phi(c, b, N)$  and  $Q(c, b, N)$ . In effect, these quantities differ by approximately constant multiplicative factors, and these factors are absorbed by  $g_2$ . Simulation results have shown that an approximation of  $L(c, b, N)$  by  $\exp(-NbH(c))$  can be good for large values of  $Nb$ . However, as mentioned above, the approximation is good only if  $b$ , the buffer per source, is quite large. For realistic specifications of an ATM switch,  $b$  may be very small, and the approximation can overestimate  $L(c, b, N)$  by several orders of magnitude. Equivalently, it overestimates the effective bandwidths.

The large  $N$  asymptotic result is expressed in the following theorem. This result has also recently been proved independently by Duffield (1996) and by Simonian and Guilbert (1994) for the case of on-off Markov fluid sources.

*Theorem 1. For the model above and under appropriate assumptions,*

$$\lim_{N \rightarrow \infty} \frac{1}{N} \log \Phi(c, b, N) = -I(c, b),$$

where

$$(5) \quad I(c, b) = \inf_m \sup_{\theta} [\theta(b + mc) - m\varphi_m(\theta)].$$

For a fixed  $m$  the supremum over  $\theta$  in (5) calculates a large deviation coefficient (appropriate to an asymptotic in  $N$ ) for the probability that the buffer overflows in  $m$  periods; the way in which this overflow occurs is by each of the  $N$  sources producing

$mc + b$  cells, and so contributing an equal share to the total of  $mC + B$  cells that causes the buffer to fill. The infimum fixes upon the  $m$  for which this probability is greatest.

*Proof of Theorem 1.* Recall that the switch is modelled as a queue with a constant service rate of  $Nc$ . If the buffer is overflowing at epoch 0 then there must have been some epoch,  $-m$ , at which the buffer was last empty and since which at least  $N(mc + b)$  cells have been received. Thus  $\Phi(c, b, N) \leq P(S_m > Nmc + Nb \text{ for some } m)$ , where  $S_m = \sum_{n=1}^m X_n$ . Pick a  $\theta_1$  such that  $c\theta_1 - \varphi(\theta_1) > 2\varepsilon$ . Then since  $\varphi_m(\theta_1) \rightarrow \varphi(\theta_1)$ , there must be an  $m_1$  such that for all  $m > m_1$ , both  $c\theta_1 - \varphi_m(\theta_1) > \varepsilon$  and  $\theta_1 b + m\varepsilon > \sup_{\theta} [\theta c - \varphi_1(\theta)]$ . A Chernoff bound is that for all  $\theta > 0$ ,

$$P(S_m > Na) \leq E \exp(\theta(S_m - Na)) = \exp(-N[\theta a - \varphi_m(\theta)]).$$

Using this, with  $a = mc + b$ , we have

$$\begin{aligned} &P(S_m > mNc + Nb \text{ for some } m) \\ &\leq \sum_{m=1}^{\infty} P(S_m > N(b + mc)) \\ &\leq \sum_{m=1}^{m_1-1} \exp \left\{ -N \sup_{\theta} [\theta(b + mc) - m\varphi_m(\theta)] \right\} + \sum_{m=m_1}^{\infty} \exp \{ -N[\theta_1 b + m\{\theta_1 c - \varphi_m(\theta_1)\}] \} \\ &\leq \sum_{m=1}^{m_1-1} \exp \left\{ -N \sup_{\theta} [\theta(b + mc) - m\varphi_m(\theta)] \right\} + \sum_{m=m_1}^{\infty} \exp \{ -N[\theta_1 b + m\varepsilon] \} \\ &\leq (m_1 - 1) \exp \left\{ -N \min_{m < m_1} \sup_{\theta} [\theta(b + mc) - m\varphi_m(\theta)] \right\} + \frac{\exp \{ -N[\theta_1 b + m_1 \varepsilon] \}}{1 - e^{-N\varepsilon}}. \end{aligned}$$

Hence

$$\begin{aligned} \overline{\lim}_{N \rightarrow \infty} \frac{1}{N} \log \Phi(c, b, N) &\leq - \min_{m < m_1} \sup_{\theta} [\theta(b + mc) - m\varphi_m(\theta)] \\ &\leq - \inf_m \sup_{\theta} [\theta(b + mc) - m\varphi_m(\theta)]. \end{aligned}$$

In the reverse direction, let us consider  $m$  consecutive epochs. The expected number of these epochs in which overflow occurs is at least as large as the probability that there is overflow during at least one of the epochs. Thus we have that, for all  $m$ ,

$$\Phi(c, b, N) \geq \frac{P(S_m > mNc + Nb)}{m}.$$

Therefore, using Cramer’s lower bound for the sum of  $N$  i.i.d. random variables

$$\underline{\lim}_{N \rightarrow \infty} \frac{1}{N} \log \Phi(c, b, N) \geq - \sup_{\theta} [\theta(b + mc) - m\varphi_m(\theta)].$$

As this holds for all  $m$ , the proof is complete.

By observing from (5) that for any given  $\delta$ ,  $NI(c, b) \geq Nb\delta$  if and only if for each  $m$  there exists a  $\theta$  such that  $\theta(b + mc) - m\varphi_m(\theta) \geq \delta b$ , we have the following equivalent expression for bandwidth requirement.

Corollary 1.1.

$$(6) \quad I(c, b) \geq b\delta \Leftrightarrow c \geq \sup_m \inf_{\theta} \left[ \frac{b}{m} \left( \frac{\delta}{\theta} - 1 \right) + \frac{\varphi_m(\theta)}{\theta} \right].$$

The far right-hand side of Equation (6) specifies the least value that may be taken by  $c$  if the large deviation coefficient  $I$  is to exceed  $b\delta$  (a number related to maximum allowable cell loss rate). Alternatively, for fixed values of  $c$  and  $\delta$ , it expresses in terms of the quantity  $\varphi_m(\theta)$  the constraint on the traffic that may be carried by the switch. Notice that there is no notion of effective bandwidth similar to that which occurs in Equation (4) for the large  $b$  asymptotic.

The large  $N$  asymptotic may be interpreted in another way: it is as if in a system, indexed by  $N$ , a single source is replaced by the average of  $N$  identically distributed sources. The asymptotic is taken as  $N \rightarrow \infty$  while the total buffer and total bandwidth are held constant. This asymptotic has been considered by Weiss (1983), for identical sources that are Markov modulated fluids. His approach is more refined in that he uses a large deviation result for sample paths to find the most likely path by which a buffer will fill. We consider now some issues raised by the large  $N$  asymptotic.

2.1. *Discretization of time.* We have assumed a discrete time model and it is illuminating to consider how this affects the results. Suppose the length of an epoch is doubled. Then  $c$  must be replaced by  $2c$ , and  $\varphi_m(\theta)$  by  $\varphi_{2m}(\theta)$ . We get  $I(c, b) = \inf_m \sup_{\theta} [\theta(b + m2c) - 2m\varphi_{2m}(\theta)]$ . So it is as if we are now taking  $\inf_m$  only over the even values of  $m$ . The answer will be greater. This is to be expected, since by taking an epoch to be longer there is some averaging of the traffic within each epoch. This dependence on the discretization is not seen in the analysis of the large  $b$  asymptotic.

2.2. *The case of no buffer.* If there is no buffer then we have the following.

Theorem 2.

$$I(c, 0) = \sup_{\theta} [\theta c - \varphi_1(\theta)].$$

*Proof.* We have

$$(7) \quad I(c, 0) = \inf_m \sup_{\theta} [\theta mc - m\varphi_m(\theta)]$$

and

$$\begin{aligned} m\varphi_m(\theta) &= \log E \left[ \exp \left( m\theta \frac{1}{m} \sum_{n=1}^m X_n \right) \right] \\ &\leq \log E \left[ \frac{1}{m} \sum_{n=1}^m \exp(m\theta X_n) \right] \\ &= \varphi_1(m\theta), \end{aligned}$$

where the first inequality follows by convexity. Thus

$$(8) \quad I(c, 0) \geq \inf_m \sup_{\theta} [\theta mc - \varphi_1(m\theta)] = \sup_{\theta'} [\theta' c - \varphi_1(\theta')].$$

The final equality follows by letting  $\theta' = m\theta$ . On the other hand, by choosing  $m = 1$  in the minimization of (7), the right-hand side in (8) is attained. Thus  $I(c, 0) = \sup_{\theta} [\theta c - \varphi_1(\theta)]$ .

Note that (7) is simply the large deviation asymptotic for the probability that the buffer should receive more than  $C$  cells in a single epoch.

2.3. *Large buffer asymptotics.* The following theorem connects the asymptotics for large  $b$  and large  $N$ .

*Theorem 3.*

$$\lim_{b \rightarrow \infty} \frac{1}{b} I(c, b) = H(c).$$

*Proof.* Recall that  $\delta^* = H(c)$  satisfies  $\varphi(\delta^*)/\delta^* = c$ . Let us take  $b_m = m(\varphi'_m(\delta^*) - c)$ . Note that since  $\varphi'_m(\delta^*) \rightarrow \varphi'(\delta^*)$ , as  $m \rightarrow \infty$ , we also have  $b_m \rightarrow \infty$  as  $m \rightarrow \infty$ . Furthermore, the growth in  $b_m$  is approximately linear for large  $m$ . Thus for any  $m$  sufficiently large that  $b_m \geq b$ ,

$$\begin{aligned} \frac{1}{b} I(c, b) &\leq \frac{b_m}{b} \frac{1}{b_m} I(c, b_m) \\ &\leq \frac{b_m}{b} \sup_{\theta} [\theta + (m/b_m)(c\theta - \varphi_m(\theta))] \\ &\leq \frac{b_m}{b} \sup_{\theta} [\theta + (m/b_m)(c\theta - \varphi_m(\delta^*) - (\theta - \delta^*)\varphi'_m(\delta^*))] \\ &= \frac{b_m}{b} \delta^* \frac{\varphi'_m(\delta^*) - \varphi_m(\delta^*)/\delta^*}{\varphi'_m(\delta^*) - \varphi(\delta^*)/\delta^*}, \end{aligned}$$

where the third inequality follows by convexity of  $\varphi_m(\theta)$ . But by the fact that  $b_m$  grows linearly in  $m$  for large  $m$ , we can let  $m \rightarrow \infty$  as  $b \rightarrow \infty$  in such a manner that  $b_m/b \rightarrow 1$ . So from the above it follows that

$$\overline{\lim}_{b \rightarrow \infty} \frac{1}{b} I(c, b) \leq \delta^*.$$

To establish an inequality in the reverse direction, pick  $\delta < \delta^*$ . Recall that  $\varphi(\delta)/\delta$  is increasing in  $\delta$  and  $\lim_{m \rightarrow \infty} \varphi_m(\delta) = \varphi(\delta)$ . Then there exists  $m_0$  such that for all  $m \geq m_0$ ,  $c\delta - \varphi_m(\delta) > 0$ . So for all  $m \geq m_0$  we have  $[b\delta + m(c\delta - \varphi_m(\delta))] > b\delta$  and

$$\begin{aligned} I(c, b) &= \inf_m \sup_{\theta} [\theta(mc + b) - m\varphi_m(\theta)] \\ &\geq \inf_m [b\delta + m(c\delta - \varphi_m(\delta))] \end{aligned}$$

$$\begin{aligned} &\geq \min \left\{ \min_{m < m_0} [b\delta + m(c\delta - \varphi_m(\delta))], b\delta \right\} \\ &= b\delta + \min_{m < m_0} [m(c\delta - \varphi_m(\delta))^-]. \end{aligned}$$

The second term on the right-hand side of the final expression is a constant and hence

$$\lim_{b \rightarrow \infty} \frac{1}{b} I(c, b) \geq \delta.$$

Since this holds for all  $\delta < \delta^*$ , the proof is complete.

Notice that in the case that  $\{X_n\}$  is a sequence of i.i.d. variables,  $\varphi_m(\theta) = \varphi(\theta)$  and so we can take  $m_0 = 0$ . Thus  $I(c, b) \geq bH(c)$  for all  $b$ . Moreover, for those  $b$  for which  $b = m(\varphi'(H(c)) - c)$ ,  $m = 1, 2, \dots$ , we have  $I(c, b) = bH(c)$ . Also, from the first part of the proof we have an upper bound, and so in the i.i.d. case,

$$H(c) \leq \frac{I(c, b)}{b} \leq \frac{b_m}{b} \delta^* = \frac{m(\varphi'(H(c)) - c)}{b} H(c) \leq \left[ 1 + \frac{\varphi'(H(c)) - c}{b} \right] H(c).$$

The following corollaries are worth noting.

*Corollary 3.1.* As  $b \rightarrow \infty$  the value of  $m$  that is optimal on the right-hand side of (5) also tends to infinity.

*Proof.* Suppose the result is not true and that, for arbitrarily large  $b$ , the optimal  $m$  is less than some  $m_0$ . Choose any  $\delta$ . Then for arbitrarily large values of  $b$ ,

$$I(c, b) > b \inf_{m < m_0} [\delta + (m/b)(c\delta - \varphi_m(\delta))].$$

This implies  $\lim_{b \rightarrow \infty} (1/b)I(c, b) > \delta$ . For  $\delta > H(c)$  this contradicts Theorem 2.

*Corollary 3.2.*

$$g_1(c, b, N)/b \rightarrow 0 \quad \text{as } b \rightarrow \infty.$$

*Proof.* This follows from  $g_1(c, b, N) = \log \Phi(c, b, N) + NI(c, b)$ . For fixed  $N$  we have that as  $b \rightarrow \infty$  both  $\log \Phi(c, b, N)/b \rightarrow -H(c)$  by de Veciana and Walrand (1994), and  $I(c, b)/b \rightarrow H(c)$  by Theorem 3. Thus  $g_1(c, b, N)/b \rightarrow 0$ .

**2.4. The shape of the acceptance region.** Recall that there are  $N\rho_i$  sources of type  $i$ . Thus the traffic is defined by the parameters  $N$  and  $\rho = (\rho_1, \dots, \rho_M)$ . Given a desired loss probability of less than  $\exp(-\delta Nb)$ , we might define as acceptable those possible mixes of sources  $\rho \in R$ , where  $R = \{\rho : Nb\delta < NI(c, b, \rho)\}$ , or equivalently,



$$(9) \quad c \geq \sup_m \inf_{\theta} \left[ \frac{b}{m} \left( \frac{\delta}{\theta} - 1 \right) + \sum_{i=1}^M \rho_i \frac{\varphi_m^i(\theta)}{\theta} \right].$$

Now we can see that in the case that the cells produced in each period are i.i.d. random variables, then modulo discreteness effects, the acceptance region is linear. For in this case,  $\varphi_m(\theta) = \varphi(\theta)$ , and

$$\begin{aligned} I(c, b) &= \inf_m \sup_{\theta} [\theta(b + mc) - m\varphi(\theta)] \\ &\geq \inf_{t>0} \sup_{\theta} [\theta(b + tc) - t\varphi(\theta)] \\ &= \sup[\theta b : \varphi(\theta)/\theta \leq c] \\ &= \sup \left[ \theta b : \sum_i \rho_i \varphi^i(\theta)/\theta \leq c \right]. \end{aligned}$$

So  $I(c, b) \geq \delta b$  if and only if  $\varphi(\delta)/\delta \leq c$ ; this is the same as (4) that holds for the large  $b$  asymptotic and shows that in this case the boundary to the acceptance region is linear in the  $\rho_i$ .

However, if we do not disregard discreteness effects, or the cells produced in successive periods are not i.i.d., the shape of the acceptance region can be more complicated. To show the effect introduced by the discreteness of  $m$  we can consider two sources in which the number of cells in successive epochs are i.i.d. Gaussian variables, with means  $\mu_1 = 5/3$ ,  $\mu_2 = 1$  and variances  $\delta_1^2 = 1/3$ ,  $\delta_2^2 = 1$ . For a source of type  $i$ ,  $\varphi^i(\theta) = \theta\mu_i + \theta^2\sigma_i^2/2$ . Hence the acceptance region defined in (9) becomes

$$c \geq \sum_{i=1}^M \rho_i \mu_i + \sup_m \left[ \sqrt{\frac{2\delta b \sum_i \rho_i \sigma_i^2}{m}} - \frac{b}{m} \right].$$

The first term is linear in  $\rho$ , but the term within the supremum is concave in  $\rho$ . For the values above, and  $\delta = 2$ ,  $b = 1$ ,  $c = 2$ , some values of  $\rho_1$ ,  $\rho_2$  that lie on the boundary of the acceptance region are shown in Table 1.

TABLE 1

$\rho_1$	0.000	0.125	0.250	0.375	0.500	0.625	0.750	0.875	1.000
$\rho_2$	1.000	0.876	0.754	0.633	0.506	0.377	0.250	0.127	0.000

These show that the boundary of the acceptance region is neither locally concave nor convex.

We believe that it is not only because of discreteness of  $m$  that the shape of the acceptance boundary is not always either concave or convex. However, we have yet to give an example of this.

Of course we know that  $I(b, c)/b \rightarrow H(c)$ , as  $b \rightarrow \infty$  and so the acceptance region defined as  $Nb\delta \geq NI$  has a boundary that is asymptotically the linear one defined in Equation (4) by  $c \geq \sum_i \rho_i \varphi^i(\delta)/\delta$ .

### 3. Markov modulated fluid sources

The Markov modulated fluid model of a source is one in which the rate of the source alternates between 0 and a peak rate  $a$  according to the state of a continuous time Markov process. The off and on states have exponentially distributed holding times with parameters  $\lambda$  and  $\mu$  respectively. The proportion of time the source is on is  $\lambda/(\lambda + \mu)$ , and so for sensible scenarios, we require  $\lambda a/(\lambda + \mu) < c < a$ .

Theorem 1 was proved under the assumption of a discrete time model, but it also clearly holds for the continuous-time setting here (cf. Duffield (1996) for a proof). To find  $I(c, b)$  for the Markov modulated fluid, we will compute the moment generating functions

$$z_{i,t}(\theta) = E_i \left[ \exp \left( \theta \int_0^t x(s) ds \right) \right], \quad i = 0, 1,$$

where  $x(t)$  is the fluid rate at time  $t$  and  $i = 0, 1$  as the source is off or on initially. Now

$$\begin{aligned} z_{0,t+\delta} &= (1 - \lambda\delta)z_{0,t} + \lambda\delta z_{1,t} + o(\delta) \\ z_{1,t+\delta} &= e^{\theta a\delta} [\mu\delta z_{0,t} + (1 - \mu\delta)z_{1,t}] + o(\delta). \end{aligned}$$

So

$$\begin{pmatrix} \dot{z}_0 \\ \dot{z}_1 \end{pmatrix} = \begin{pmatrix} -\lambda & \lambda \\ \mu & \theta a - \mu \end{pmatrix} \begin{pmatrix} z_0 \\ z_1 \end{pmatrix}$$

and

$$\begin{aligned} z_{0,t} &= \frac{\omega_2}{\omega_2 - \omega_1} e^{\omega_1 t} - \frac{\omega_1}{\omega_2 - \omega_1} e^{\omega_2 t} \\ z_{1,t} &= \frac{\omega_2 - \theta a}{\omega_2 - \omega_1} e^{\omega_1 t} + \frac{\theta a - \omega_1}{\omega_2 - \omega_1} e^{\omega_2 t}, \end{aligned}$$

where

$$\begin{aligned} \omega_1 &= \frac{-(\lambda + \mu - \theta a) - \sqrt{(\lambda + \mu - \theta a)^2 + 4\lambda\theta a}}{2} \\ \omega_2 &= \frac{-(\lambda + \mu - \theta a) + \sqrt{(\lambda + \mu - \theta a)^2 + 4\lambda\theta a}}{2} \end{aligned}$$

are the two roots of  $\omega^2 + (\lambda + \mu - \lambda\theta a)\omega - \lambda\theta a = 0$ . Hence

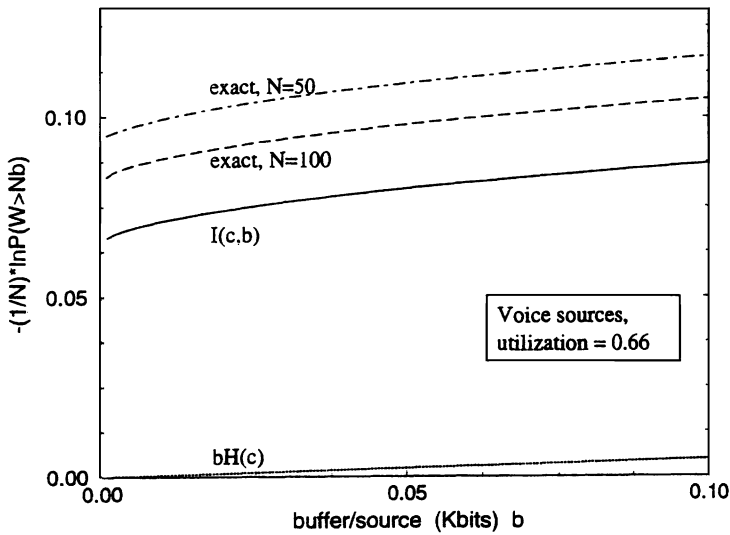


Figure 1. Markov modulated fluid sources

$$\begin{aligned}
 z_t(\theta) &= E \left[ \exp \left( \theta \int_0^t x(s) ds \right) \right] \\
 &= \frac{\mu}{\lambda + \mu} z_{0,t}(\theta) + \frac{\lambda}{\lambda + \mu} z_{1,t}(\theta) \\
 &= \frac{\mu\omega_2 + \lambda(\omega_2 - \theta a)}{(\omega_2 - \omega_1)(\lambda + \mu)} e^{\omega_1 t} + \frac{-\mu\omega_1 + \lambda(\theta a - \omega_1)}{(\omega_2 - \omega_1)(\lambda + \mu)} e^{\omega_2 t}.
 \end{aligned}$$

The problem which we wish to solve is  $I(c, b) = \inf_{t>0} \sup_{\theta} [\theta(b + tc) - \log(z_t(\theta))]$ .

By Theorem 3, we know that  $I(c, b)/b \rightarrow H(c)$ , and also that the  $t$  that is optimal in the right-hand side above tends to  $\infty$  as  $b \rightarrow \infty$ . Note that  $\varphi(\theta) = \lim_{t \rightarrow \infty} t^{-1} \log z_t(\theta) = \omega_2(\theta)$ . The bandwidth requirement  $c \geq \omega_2(\delta)/\delta \Leftrightarrow H(c) \geq \delta$  is precisely that given in other papers, such as Gibbens and Hunt (1991). An alternative expression of this constraint is

$$H(c) = \frac{(\lambda + \mu)c - \lambda a}{c(a - c)} \leq \delta.$$

3.1. *Numerical results.* We have conducted a number of experiments to determine the effectiveness of estimating buffer overflow rates using the large  $N$  asymptotic. For a Markov modulated fluid model,  $Q(c, b, N) = P(W > Nb)$  can be obtained exactly using techniques of Anick *et al.* (1982). To model a voice source we use the parameters  $\lambda = 650\text{ms}$ ,  $\mu = 353\text{ms}$  and  $a = 64\text{kbps}$ . The mean bandwidth is 22.48kbps. In Figure 1, we take  $c$  to be as 1.5 times the mean rate, i.e. 33.72kbps. We plot as functions of  $b$  graphs of  $I(c, b)$ ,  $bH(c)$  and the exact value of  $(1/N)\log Q(c, b, N)$  as computed using the

formula of Anick *et al.* (1982) for  $N = 50, 100$ . Notice that the large  $N$  asymptotic is a much better estimator of  $Q(c, b, N)$  than is the large  $b$  asymptotic, and that it is more exact for larger  $N$ . For example, at  $b = 0.05$ ,  $-(1/N)\log Q(c, b, N)$  is 0.10903 for  $N = 50$ , 0.09758 for  $N = 100$ ; these are well-estimated by  $I(c, b) = 0.07989$ , but not by  $bH(c) = 0.00241$ . Both estimates are conservative in the sense that they both overestimate the true probability that the queue should exceed  $Nb$ .

3.2. *The shape of  $I(c, b)$  near  $b = 0$ .* Our numerical results agree with the calculations of Weiss (1983) that for the Markov modulated fluid the behavior for small  $b$  is  $I(c, b) = C_1 + C_2\sqrt{b}$ . We have not been able to establish the  $\sqrt{b}$  behavior from (5), but we can compute the constant  $C_1 = I(c, 0)$ . Following the ideas in Theorem 2, suppose  $s$  is very large. Then

$$\begin{aligned} I(c, 0) &= \inf_{s>t>0} \sup_{\theta} [\theta tc - t\varphi_t(\theta)] \\ &\geq \inf_{s>t>0} \sup_{\theta} [n\theta(t/n)c - (t/n)\varphi_{t/n}(n\theta)] \\ &= \inf_{s/n>t>0} \sup_{\theta} [\theta tc - t\varphi_t(\theta)] \\ &\geq \lim_{t \rightarrow 0} \sup_{\theta} [\theta tc - t\varphi_t(\theta)] \\ &= \lim_{t \rightarrow 0} \sup_{\theta} [\theta tc - \log E[\exp(\theta[tX(0) + o(t)])]] \\ &= \lim_{t \rightarrow 0} \sup_{\theta} \left[ \theta c - \log \left( \frac{\lambda e^{\theta b}}{\lambda + \mu} + \frac{\mu}{\lambda + \mu} + o(t)/t \right) \right] \\ &= c' \log \left( \frac{c'}{p} \right) + (1 - c') \log \left( \frac{1 - c'}{1 - p} \right), \end{aligned}$$

where  $p = \lambda/(\lambda + \mu)$  and  $c' = c/a$ . Clearly  $I(c, 0) \leq \lim_{t \rightarrow 0} \sup_{\theta} [\theta tc - t\varphi_t(\theta)]$ . So this is in fact the value of  $I(c, 0)$ , and it is just the asymptotic rate for the probability that the number of on sources exceeds  $c/a$ .

#### 4. Gaussian sources

4.1. *Gaussian stationary sources.* Suppose  $\{X_n\}$  is the superposition of  $N$  Gaussian sources, each distributed as  $\{Y_n\}$  with mean  $\mu$ , variance  $\sigma^2$  and autocovariance function  $\gamma(k)$ . Then  $\varphi_m(\theta) = \mu\theta + \sigma_m^2\theta^2/2$ , where  $m\sigma_m^2$  is given by

$$m\sigma_m^2 = \text{var} \left( \sum_{i=1}^m Y_i \right) = m\sigma^2 + 2[(m-1)\gamma(1) + (m-2)\gamma(2) + \dots + \gamma(m-1)].$$

Notice that  $\lim_{m \rightarrow \infty} \sigma_m^2 = \gamma$ , where  $\gamma$  is the index of dispersion and equal to  $\sum_{-\infty}^{\infty} \gamma(k)$ .  
Now

$$I(c, b) = \inf_m \frac{(b + mc - m\mu)^2}{2m\sigma_m^2},$$

with the optimum achieved by  $\theta = (b + mc - m\mu)/m\sigma_m^2$ .

Also the bandwidth requirement (6) is

$$c \geq \sup_m \inf_\theta \left[ \frac{b}{m} \left( \frac{\delta}{\theta} - 1 \right) + \frac{\varphi_m(\theta)}{\theta} \right] = \sup_m \left[ \mu + \sqrt{\frac{2b\delta\sigma_m^2}{m}} - \frac{b}{m} \right],$$

with the optimum achieved by  $\theta = \sqrt{2b\delta/m\sigma_m^2}$ .

4.2. *An autoregressive source.* A simple autoregression has sometimes been used as a model for a videotelephone source. Let us take  $Y_n = \alpha Y_{n-1} + (1 - \alpha)\mu + \varepsilon_n$ , where  $\{\varepsilon_n\}$  is Gaussian white noise with variance  $\eta^2$ . Then  $\sigma^2 = \eta^2/(1 - \alpha^2)$ ,  $\gamma = (1 + \alpha)\sigma^2/(1 - \alpha)$  and  $\sigma_m^2$  can be derived as follows:

$$\begin{aligned} m\sigma_m^2 &= (m - 1)\sigma_{m-1}^2 + \sigma^2 + 2(\alpha + \alpha^2 + \dots + \alpha^{m-1})\sigma^2 \\ &= (m - 1)\sigma_{m-1}^2 + 2\sigma^2 \frac{1 - \alpha^m}{1 - \alpha} - \sigma^2 \\ &= \sigma^2 + 2\sigma^2 \frac{(1 - \alpha^m) + \dots + (1 - \alpha^2)}{1 - \alpha} - (m - 1)\sigma^2 \\ &= -(m - 2)\sigma^2 + \frac{2\sigma^2}{1 - \alpha} \left( m - 1 - \frac{\alpha^2 - \alpha^{m+1}}{1 - \alpha} \right) \\ &= m \frac{1 + \alpha}{1 - \alpha} \sigma^2 - 2 \frac{\alpha - \alpha^{m+1}}{(1 - \alpha)^2} \sigma^2 \\ &= m\gamma - 2 \frac{\alpha(1 - \alpha^m)}{(1 - \alpha)^2} \sigma^2. \end{aligned}$$

It turns out that  $\sigma_m^2$  is a convex increasing function of  $m$  when  $\alpha > 0$ . Notice that

$$\sigma_m^2 = \gamma - 2 \frac{\alpha(1 - \alpha^m)}{m(1 - \alpha)^2} \sigma^2 < \gamma$$

if  $\alpha > 0$ . But  $\sigma_m^2 > \gamma$  for  $\alpha < 0$ . Now since  $\varphi(\theta)$  is quadratic in  $\theta$ , we can easily find that

$$\sup_\theta [\theta(b + mc) - m\varphi_m(\theta)] = \frac{(b + mc - m\mu)^2}{2m\gamma - 4 \frac{\alpha(1 - \alpha^m)}{(1 - \alpha)^2} \sigma^2}.$$

For large  $b$  this is minimized by large  $m$ , and so assuming  $\alpha^m$  is small this is minimized with respect to  $m$  by

$$m \simeq \frac{b}{c - \mu} + \xi / \gamma,$$

where

$$\xi = 4 \frac{\alpha}{(1 - \alpha)^2} \sigma^2.$$

The optimal value is

$$I(c, b) \simeq \frac{2b(c - \mu)}{\gamma} + \frac{\xi(c - \mu)^2}{\gamma^2} = \frac{2b(c - \mu)}{\gamma} + \frac{4\alpha(c - \mu)^2}{(1 + \alpha)^2 \sigma^2}.$$

This can be compared with  $bH(c) = 2b(c - \mu)/\gamma$ . One sees that  $bH(c) < I(c, b)$  when  $\alpha > 0$ , but that the inequality is reversed when  $\alpha < 0$ . The case  $\alpha > 0$  corresponds to a relatively bursty source, while  $\alpha < 0$  corresponds to relatively smooth source, because of the negative first order autocorrelation. When  $\alpha > 0$  the use of the large  $b$  asymptotic to estimate  $L(c, b, N)$  leads to a larger estimate of cell loss rate than does the large  $N$  asymptotic: the experimental evidence is that both asymptotics lead to overestimates of the true cell loss rate but that the large  $N$  asymptotic is closer. When  $\alpha < 0$  the large  $b$  asymptotic estimates a smaller loss rate than the estimate based on a large  $N$  asymptotic: the experimental evidence is that the large  $b$  asymptotic underestimates the true loss rate while the large  $N$  asymptotic overestimates it and is closer. In this case the large  $N$  asymptotic is clearly preferable.

4.3. *Numerical results.* For an autoregressive Gaussian process with  $\alpha = -0.5$ ,  $\mu = 18$ ,  $\eta^2 = 64$  and  $c = 18.15$ , we plot  $\exp(-NbI(b, c))$  and  $\exp(-bH(c))$  against  $b$ , together with simulation results (over 2,000,000 epochs) for the probability  $Q(c, b, N) = P(W > Nb)$ , for  $N = 1000$  and five values of  $b$ , 1.0, 1.2, 1.4, 1.6, 1.8. The proportion of the time that the buffer is empty is 60.95%. Figure 2 shows that for these small amounts of buffer per source, the large  $N$  asymptotic overestimates the tail probability by a factor of about 2–3, while the large  $N$  asymptotic underestimates it by a factor of about 2–5. For example, at  $b = 1.2$ , the true value has a 95% confidence interval  $(92, 104) \times 10^{-5}$ , and mean  $98 \times 10^{-5}$ ; the large  $N$  asymptotic gives  $263 \times 10^{-5}$  and the large  $b$  asymptotic gives  $32 \times 10^{-5}$ .

4.4. *The sign of  $I(c, b) - bH(c)$ .* The above example is illuminating in considering the difference between the large  $N$  and large  $b$  asymptotics. Our results are consistent with the work of Choudhury *et al.* (1994a), (1994b), who discuss the fact that their numerical experience suggests that the large  $B$  asymptotic should be modified by prior multiplication by an exponential factor in  $N$ , so  $Q(c, b, N) \sim \beta e^{-N\gamma} e^{-\eta B}$ . This can now be explained theoretically by writing

$$\begin{aligned} Q(c, b, N) &= \exp\{-NI(c, b) + g_1(c, b, N)\} \\ &= \exp\{-N[I(c, b) - bH(c)] - NbH(c) + g_1(c, b, N)\} \\ &= \exp\{-N[I(c, b) - bH(c) - g_1(c, b, N)/N]\} \exp(-BH(c)). \end{aligned}$$

Of course  $\eta = H(c)$  and  $g_1(c, b, N)/N \rightarrow 0$  as  $N \rightarrow \infty$  with  $b$  fixed.

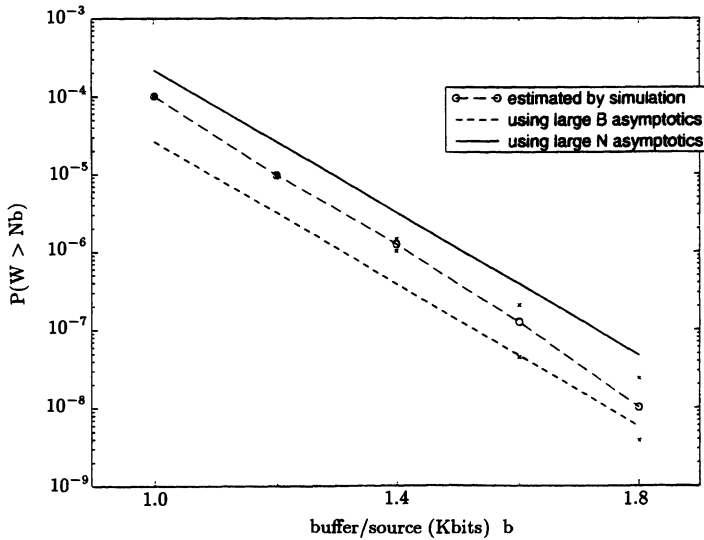


Figure 2. Autoregressive Gaussian sources

It is therefore clear that  $\gamma$  has the sign of  $I(c, b) - bH(c) - g_1(c, b, N)/N$  which for fixed  $b$  is determined by the sign of  $I(c, b) - bH(c)$  for large  $N$ . We have seen in the discussion above that  $I(c, b) > bH(c)$  when  $\alpha > 0$  and  $I(c, b) < bH(c)$  when  $\alpha < 0$ , and that these correspond to greater or less burstiness of the source, and give  $\gamma > 0$  and  $\gamma < 0$  respectively. Notice that  $I(c, b) > bH(c)$  if  $\varphi_m(\theta) > \varphi(\theta)$ , for all  $m$  and  $\theta$ . Note that if for two random variables,  $X$  and  $Y$ ,  $E \exp(\theta X) > E \exp(\theta Y)$  then  $X$  is more variable than  $Y$ . Similarly, an ordering of moment generating functions corresponds to an ordering of bustiness. This means that the quantity of cells arriving over a number of periods  $m$  is actually less bursty than is implicitly assumed when one uses  $\varphi(\theta)$  in (3) to calculate  $H(c)$ . Therefore  $\exp(-bH(c))$  tends to overestimate  $Q(c, b, N)$ .

In the cases we have considered it appears that  $\varphi_m(\theta)$  is monotone in  $m$ , and thus  $\gamma > 0$  and  $\varphi_m(\theta) < \varphi(\theta)$  corresponds to the case in which  $\varphi_1(\theta) < \varphi_m(\theta)$ , and in which we would say that the source is ‘more bursty than Poisson’, in the sense that for a Poisson source, or any which is i.i.d., we have  $\varphi_1(\theta) = \varphi_m(\theta)$ . The same comments apply, with inequalities reversed, so  $\gamma < 0$  occurs for sources that are ‘less bursty than Poisson’.

### 5. Open issues and remarks

5.1. *On-line estimation.* An important issue in ATM networks is the accurate and timely estimation of the cell-loss probabilities that occur at various switches in the network, so that correct decisions can be made about whether or not to accept more traffic or how it should be routed. A generic problem that the network management system faces is that because the events of interest are rare, e.g. cell loss rates of the order of  $10^{-6}$ – $10^{-10}$ , any brute-force on-line estimation procedure would fail because of the large time that would be required to make an accurate estimate. Providing a reasonable

confidence interval for the estimator requires time of the order of magnitude of tens of minutes, which is impractical since decisions about accepting new traffic in the network must be done immediately, and because the composition of the traffic can change in such a large time interval.

The large  $N$  asymptotic can be used to estimate in real-time the cell loss rate in the following manner. Suppose  $N\rho$  sources (comprising  $N\rho_i$  sources of type  $i$ ) are active in a switch with total bandwidth  $C$  and total buffer  $B$ . Suppose that we can make an on-line measurement to estimate the cell loss rate that would occur if  $N\rho/k$  sources were routed through a switch that has total bandwidth  $C/k$  and total buffer  $B/k$ . This can be done by a special device at the switch level which simulates a 'virtual' switch of  $1/k$  the size of the actual one, operating on a representative sample of the actual input of size  $1/k$ . The device provides an estimator of the buffer overflow rate that occurs in the virtual system. The large  $N$  asymptotic suggests that the loss rate in the virtual system should be about  $p = \exp(-NI(c, b)/k)$ , hence it can be estimated accurately in a time that is orders of magnitude smaller than the time required to make an estimate of the buffer overflow rate in the original system, and we can extrapolate the actual loss rate to be  $p^k$ .

For example, suppose the channel is carrying 1000 sources, with a loss rate of about  $10^{-10}$ . As we have already mentioned, such a small cell loss rate is difficult to measure directly. However, an on-line measurement of the loss rate for 200 sources in a virtual system with one-fifth the bandwidth and buffer space should measure a loss rate  $p$  of about  $10^{-2}$ . At this loss rate cell losses will be observed in a relatively short time and the loss rate can be measured satisfactorily. Then  $p^5$  is an estimate of the actual cell loss rate. Of course there is the opportunity to make five independent estimates of  $p$ , using five groups of 200 sources.

This idea can be refined further to provide a more accurate estimate of the cell-loss probability by computing the  $o(N)$  terms in (1). One can assume that  $\Phi(N) = AN^\psi \exp(-NI)$  and estimate the values of  $A$ ,  $\psi$ ,  $I$  by measuring  $\Phi$  in three virtual systems of size  $N/k_1$ ,  $N/k_2$ ,  $N/k_3$ , where  $N/k_1 < N/k_2 < N/k_3 \ll N$ . Then  $\Phi(N)$  is computed by substituting the corresponding values. This provides an alternative to the MINOS procedure described in Courcoubetis *et al.* (1995). Numerical work has validated the utility of this heuristic approach.

**5.2. Traffic of random composition.** The class of traffic may not be known exactly. Suppose that in a certain class of traffic the sources are actually of  $k$  possible types, with each source having probability  $p_i$  of being of type  $i$ . The probability of buffer overflow is maximized when the observed mix deviates from the expected mix in an appropriate way and sources produce cells at above their mean rates. The probability that  $N$  sources of such random mixture should have empirical distribution  $\pi_1, \dots, \pi_k$  is

$$P(\pi) = \exp\left(-N \sum_{i=1}^k \pi_i \log(\pi_i/p_i) + o(N)\right),$$

and thus taking the product of probabilities,



$$\lim_{N \rightarrow \infty} \frac{1}{N} \log \Phi(c, b) = -\inf_m \inf_{\pi_i: \sum_{i=1}^k \pi_i = 1} \sup_{\theta} \left[ \theta(b + mc) - \sum_{i=1}^k \pi_i \{ \varphi_m^i(\theta) - \log(\pi_i/p_i) \} \right].$$

Assuming the optimization over  $\pi_i$  occurs at an interior point, the stationarity condition implies  $\pi_i$  is proportional to  $p_i \exp(\varphi_m^i(\theta))$ , and hence after some algebra,

$$\begin{aligned} I_{\pi}(c, b) &= \lim_{N \rightarrow \infty} \frac{1}{N} \log \Phi(c, b) \\ (10) \qquad &= -\inf_m \sup_{\theta} \left[ \theta(b + mc) - \log \left( \sum_{i=1}^k p_i E[\exp(\theta S_m^i)] \right) \right] \\ &= -\inf_m \sup_{\theta} \left[ \theta(b + mc) - \log \left( \sum_{i=1}^k p_i \exp(m\varphi_m^i(\theta)) \right) \right], \end{aligned}$$

where  $S_m^i$  has the stationary distribution of the number of cells produced by a source of type  $i$  over  $m$  epochs. Note that the final term in (10) is, by convexity,

$$\log \left( \sum_{i=1}^k p_i \exp(m\varphi_m^i(\theta)) \right) \geq \sum_{i=1}^k p_i m\varphi_m^i(\theta) = m\varphi_m(\theta),$$

and thus the asymptotic overflow frequency, given by (10), is greater than when the mix is known exactly. Then as  $b \rightarrow \infty$  the optimal value of  $m$  tends to infinity and the sum in (10) is dominated by the maximal term, hence

$$\lim_{b \rightarrow \infty} \frac{I_{\pi}(c, b)}{b} = H_{\pi}(c) = \max \{ \theta : \varphi_i(\theta)/\theta < c, \text{ for all } i \} = \min_i H_i(c),$$

where  $H_i(c)$  is obtained for a source consisting only of type  $i$ . Similarly,  $H_{\pi}(c) > \delta$  if and only if  $\varphi_i(\delta)/\delta < c$  for all  $i$ , and so this analysis concludes that the effective bandwidth of the random mix is equal to the maximum of the effective bandwidths of the individual source types. This makes sense since the most likely way that a single source of unknown type will fill a large buffer is if the source type turns out to be the most bursty of the  $k$  possible types.

*5.3. Open problems.* A number of further issues require study. These include: (a) the efficient numerical solution of the double optimization problem defining  $I(c, b)$ ; it is not clear that the function always has a saddle point, though this is the case in the examples we have studied; (b) further understanding of the properties of the solution and the implied acceptance region; and (c) the on-line estimate of the  $\varphi_m(\theta)$ , or the on-line estimation of  $I(c, b)$ .

**Acknowledgment**

We are very grateful to Georgos Fouskas for computing the numerical result reported in Sections 4 and 5.

## References

- ANICK, D., MITRA, D. AND SONDDHI, M. (1982) Stochastic theory of a data-handling system with multiple sources. *Bell System Tech. J.* **61**, 1872–1894.
- CHOU DHURY, G., LUCANTONI, D. AND WHITT, W. (1994a) Squeezing the most out of ATM. *IEEE Trans. Commun.* to appear
- CHOU DHURY, G., LUCANTONI, D. AND WHITT, W. (1994b) On the effectiveness of effective bandwidths for admission control in ATM networks. In *ITC 14*, ed. J. Labetouille and J. Roberts. Elsevier, Amsterdam. pp. 411–420.
- COURCOUBETIS, C., KESIDIS, G., RIDDER, A., WALRAND, J. AND WEBER, R. (1995) Admission control and routing in ATM networks using inferences from measured buffer occupancy. *IEEE Trans. Commun.* **43**, 1778–1784.
- COURCOUBETIS, C. AND WEBER, R. (1995) Effective bandwidths for stationary sources. *Prob. Eng. Inf. Sci.* **7**, 285–296.
- DE VECIANA, G. AND WALRAND, J. (1994) Effective bandwidths: call admission, traffic policing and filtering for ATM networks. *Queueing Systems* **20**, 37–59.
- DUFFIELD, N. (1996) Economies of scale in queues with sources having power-law large deviation scalings. *J. Appl. Prob.* **33**, 840–857.
- ELWALID, A. AND MITRA, D. (1993) Effective bandwidth of general Markovian traffic sources and admission control of high speed networks. *IEEE/ACM Trans. Network.* **1**, 329–343.
- GIBBENS, R. AND HUNT, P. (1991) Effective bandwidths for the multi-type UAS channel. *Queueing Systems* **1**, 17–28.
- KELLY, F. (1991) Effective bandwidths at multi-class queues. *Queueing Systems* **9**, 5–16.
- KESIDIS, G., WALRAND, J. AND CHANG, C. (1993) Effective bandwidths for multiclass Markov fluids and other ATM sources. *IEEE/ACM Trans. Network.* **1**, 424–428.
- SIMONIAN, A. AND GUILBERT, J. (1994) Large deviations approximation for fluid queues fed by a large number of on-off sources. Preprint.
- WEISS, A. (1983) The large deviation of a Markov process which models traffic generation. *Technical report*. AT&T Bell Laboratories.
- WHITT, W. (1993) Tail probabilities with statistical multiplexing and effective bandwidths in multi-class queues. *Telecommun. Systems* **2**, 71–107.