

# Effective Bandwidths for Multiclass Markov Fluids and Other ATM Sources

George Kesidis, *Member, IEEE*, Jean Walrand, *Fellow, IEEE*, and Cheng-Shang Chang, *Senior Member, IEEE*

**Abstract**—We show the existence of effective bandwidths for multiclass Markov fluids and other types of sources that are used to model ATM traffic. More precisely, we show that when such sources share a buffer with deterministic service rate, a constraint on the tail of the buffer occupancy distribution is a linear constraint on the number of sources. That is, for a small loss probability one can assume that each source transmits at a fixed rate called its effective bandwidth. When traffic parameters are known, effective bandwidths can be calculated and may be used to obtain a circuit-switched style call acceptance and routing algorithm for ATM networks. The important feature of the effective bandwidth of a source is that it is a characteristic of that source and the acceptable loss probability only. Thus, the effective bandwidth of a source does not depend on the number of sources sharing the buffer or the model parameters of other types of sources sharing the buffer.

## I. INTRODUCTION

EFFECTIVE bandwidths have been discovered for certain traffic models and certain performance criteria (see [17],[13],[12],[4],[18],[3]). For example, consider a buffer of infinite size with service rate  $c$  cells/second. Assume that the buffer sources and occupancy are in steady state. Let  $X$  be the number of cells in the buffer found by a typical arriving cell. Suppose that

$$P\{X > B\} \leq e^{-B\delta} \quad (1)$$

must be satisfied (the performance criterion). Suppose further that there are  $N_j$  independent on/off Markov fluids [1] of type  $j$  ( $j = 1, 2, \dots, K$ ) sharing the buffer. There exist functions  $\alpha_j$  that depend only on the parameters of a type  $j$  source and  $\delta$ , such that the constraint (1) holds for  $B\delta \gg 1$  if and only if

$$\sum_{j=1}^K N_j \alpha_j \leq c.$$

We call  $\alpha_j$  the *effective bandwidth* of an on/off Markov fluid of type  $j$  (see [12] and [13] for proofs of this result and numerical examples that explore the accuracy of the effective bandwidth approach).

In general, effective bandwidths depend on both the traffic/buffer models and the performance criterion. Kelly [17]

Manuscript received June 1992; revised October 1992 and April 1993; recommended for transfer from the IEEE TRANSACTIONS ON COMMUNICATIONS by the IEEE/ACM TRANSACTIONS ON NETWORKING Editor-in-Chief.

G. Kesidis is with the E&CE Department, University of Waterloo, Ontario, Canada. (email: kesidis@cheetah.vlsi.uwaterloo.ca)

J. Walrand is with the EECS Department, University of California at Berkeley, Berkeley, CA. (email: wlr@diva.berkeley.edu)

C.-S. Chang is with the IBM T.J. Watson Research Center, Yorktown Heights, NY. (email: cschang@watson.ibm.com)

IEEE Log Number 9211588.

finds effective bandwidths for GI/G/1 queues under (1) and for M/G/1 queues with the performance criterion taken to be the buffer utilization (fraction of time  $X \neq 0$ ) or mean workload ( $EX < B$ ). Courcoubetis and Walrand [4] find effective bandwidths for stationary Gaussian sources under (1). Recently, Elwalid and Mitra [8] obtained effective bandwidth results for the case of continuous-time Markovian sources under (1) (c.f., Sections III-C and III-D and the Conclusions). The open question answered in this note is the existence of effective bandwidths for more general source models under (1).

We start by heuristically deriving an expression for  $P\{X > B\}$  for general source models. Consider an infinite buffer with service rate  $c$  shared by  $N_i$  sources of type  $i$ ,  $i = 1, \dots, K$ . All the sources are assumed independent. For all  $M_i$  greater than the average rate of cells produced by a source of type  $i$ , assume that the probability that a source of type  $i$  produces  $M_i T$  cells over a period of time of length  $T$  is approximately  $\exp(-TH_i(M_i))$  where  $H_i$  is convex and nonnegative (this assumption is motivated by the theory of large deviations and will be discussed). This approximation is sharpest for  $T \gg 1$ .

By independence, the probability that for  $j = 1, \dots, N_i$  the  $j^{\text{th}}$  source of type  $i$  produces  $\mu_j T$  cells over time  $T$  is about

$$\exp\left(-T \sum_{j=1}^{N_i} H_i(\mu_j)\right).$$

Consequently, the probability that all sources of type  $i$  produce a total of  $N_i M_i T$  cells over large time  $T$  is about

$$\sum_{\mu: \sum \mu_j = N_i M_i} \exp\left(-T \sum_{j=1}^{N_i} H_i(\mu_j)\right)$$

where  $\mu = (\mu_1, \dots, \mu_{N_i})$ . Indeed, each choice of  $\mu$  such that  $\sum \mu_j = N_i M_i$  is one particular way for  $N_i M_i T$  cells to get produced. This sum of exponentials can be approximated by the largest term (originally an argument of Laplace):

$$\begin{aligned} & \sum_{\mu: \sum \mu_j = N_i M_i} \exp\left(-T \sum_{j=1}^{N_i} H_i(\mu_j)\right) \\ & \approx \exp\left(-\inf_{\mu: \sum \mu_j = N_i M_i} T \sum_{j=1}^{N_i} H_i(\mu_j)\right) \\ & = \exp(-TN_i H_i(M_i)) \end{aligned}$$

where the last equality is due to the convexity of  $H_i$ . Therefore, by independence, the probability that for  $i = 1, \dots, K$  the

sources of type  $i$  produce  $N_i M_i T$  cells over time  $T$  is about

$$\exp\left(-T \sum_{i=1}^K N_i H_i(M_i)\right).$$

Thus, the probability that, starting from an empty buffer, the sources of type  $i$  produce cells at rate  $N_i M_i$  until the buffer occupancy exceeds  $B$  is

$$\exp\left(-B \frac{\sum N_i H_i(M_i)}{\sum N_i M_i - c}\right).$$

Indeed,  $T = B/(\sum N_i M_i - c)$  is the time the buffer occupancy takes to reach  $B$  when the aggregate cell arrival rate is  $\sum N_i M_i$ . By the argument of Laplace, the probability that the buffer occupancy, starting from empty, reaches  $B$  before it returns to empty is about

$$\exp\left(-B \inf_{\sum N_i M_i > c} \frac{\sum N_i H_i(M_i)}{\sum N_i M_i - c}\right) \approx P\{X > B\}. \quad (2)$$

Given the effective bandwidths of a buffer's sources, one can determine its *spare capacity* to accept more calls at any time. For instance, say we want to determine if a call of type  $j$  can be accommodated (i.e., constraint (1) is preserved) in a buffer that is currently being used by  $N_i$  calls of type  $i$ ,  $i = 1, \dots, K$ . If  $\alpha_j(\delta) < c - \sum_{i=1}^K N_i \alpha_i(\delta)$ , then the call can be accommodated; otherwise, it cannot. See, for example, [10], [13], [14], and [19] for further discussion on how effective bandwidths can be used for network resource management.

This note is organized as follows. In Section II, we show the existence of effective bandwidths in the multiclass case when the sources satisfy certain conditions. In Section III, we give expressions for the effective bandwidths of Markov-modulated Poisson processes, Markov-modulated fluids (or just "Markov fluids"), and discrete-time Markov sources. Finally, conclusions are drawn in Section IV.

## II. GENERAL EFFECTIVE BANDWIDTHS

We now show the existence of effective bandwidths. First some assumptions on the sources are made, then effective bandwidths are defined by considering the single-source case, and finally the multiclass case is considered.

Consider an infinite buffer with deterministic service rate  $c$  cells/second, shared by  $N_i$  independent sources of type  $i$ ,  $i = 1, \dots, K$ . Denote by  $[\cdot, \cdot]$  the scalar product. Let  $\Gamma_i \in (0, \infty]$  (respectively  $\gamma_i \in [0, \infty)$ ) denote the maximum (respectively minimum) possible cell arrival rate of a type  $i$  source. Let  $\bar{\gamma}_i \in (0, \infty)$  be the average arrival rate of a type  $i$  source. We assume that

$$N \in \mathbf{C} := \{N \in \mathbf{Z}_+^K : [N, \Gamma] > c \text{ and } [N, \bar{\gamma}] < c\}$$

where  $\mathbf{Z}_+ = \{0, 1, 2, \dots\}$ ,  $\bar{\gamma} = (\bar{\gamma}_1, \dots, \bar{\gamma}_K)$ , and  $\Gamma := (\Gamma_1, \dots, \Gamma_K)$ . Let  $M = (M_1, \dots, M_K)$ .

Motivated by (2), we take the measure of congestion in the buffer to be

$$\exp(-BI(N, c) + o(B)) \quad (3)$$

where

$$I(N, c) := \inf_{M \in \mathbf{A}(N, c)} \frac{\sum_{i=1}^K N_i H_i(M_i)}{[N, M] - c}$$

and  $\mathbf{A}(N, c) := \{M \in \mathbf{R}_+^K : \gamma_i < M_i < \Gamma_i \forall i \text{ and } [N, M] > c\}$  (c.f., (5) for the definition of  $H_i$ ). Thus, when  $B\delta \gg 1$ , the constraint (1) is

$$I(N, c) \geq \delta. \quad (4)$$

Assume that the sources are stationary and ergodic. Consider a single source of type  $i$ . Let the number of arrivals of this type  $i$  source in the time interval  $[0, t]$  be  $A_i(t)$ . Assume that  $A_i$  satisfies the conditions of the Gartner-Ellis Theorem [11], [6], [2]. That is, assume that the asymptotic log moment generating function of  $A_i$ ,

$$h_i(\delta) := \lim_{t \rightarrow \infty} \frac{1}{t} \log \mathbf{E} \exp(A_i(t)\delta)$$

exists and is finite for all real  $\delta$ , and that  $h_i$  is differentiable. We can directly verify that  $h_i$  is convex, positive, and increasing for  $\delta > 0$ .

By the Gartner-Ellis Theorem,  $H_i$  is the Legendre transform of  $h_i$ :

$$H_i(M_i) := \sup_{\delta \in \mathbf{R}} \{\delta M_i - h_i(\delta)\}. \quad (5)$$

We can directly verify that  $H_i$  is nonnegative, convex, differentiable,  $H_i(\bar{\gamma}_i) = 0$ , and  $H_i(M) = \infty$  for all  $M > \Gamma_i$  or  $M < \gamma_i$ . We also assume that  $H_i$  is *strictly* convex on the interval  $(\gamma_i, \Gamma_i)$ .

Consider the case of a single source of type  $i$ . For  $\delta > 0$ , define  $\alpha_i(\delta)$  to be the value of  $a$  such that

$$I_i(a) := \inf_{M_i \in \mathbf{A}_i(a)} \frac{H_i(M_i)}{M_i - a} = \delta$$

where  $\mathbf{A}_i(a) := \{M_i : a < M_i < \Gamma_i\}$ . Thus,  $\alpha_i(\delta) = I_i^{-1}(\delta)$  can be interpreted as the rate at which to serve a single source of type  $i$  so that constraint (4) is satisfied. We call  $\alpha_i$  the *effective bandwidth* of type  $i$  traffic. The following theorem gives us a more manageable form for  $\alpha_i$ .

*Lemma 1:* Under these conditions, for all  $\delta > 0$ ,

$$\alpha_i(\delta) = \frac{h_i(\delta)}{\delta}.$$

*Proof:* Since  $H_i$  and  $h_i$  are convex conjugates,  $h_i(\delta) = \sup_M \{M\delta - H_i(M)\}$ . It then follows from the differentiability of  $H_i$  and  $h_i$  and the strict convexity of  $H_i$  that

$$h_i(\delta) = \delta H_i'^{-1}(\delta) - H_i(H_i'^{-1}(\delta)). \quad (6)$$

Define the function  $g_i(M) := M - H_i(M)/H_i'(M)$ . From the strict convexity of  $H_i$ , it follows that  $g_i$  is strictly increasing on  $(\gamma_i, \Gamma_i)$ . Thus, we can define  $g_i^{-1}$  as the inverse of  $g_i$ ; i.e., for  $a \in (\bar{\gamma}_i, \Gamma_i)$ ,  $g_i^{-1}(a)$  is the solution of  $a = M - H_i(M)/H_i'(M)$ . Since  $H_i' > 0$  on  $(\bar{\gamma}_i, \Gamma_i)$ ,  $g_i^{-1}(a) > a$  so that  $g_i^{-1}(a) \in \mathbf{A}_i(a)$ . Thus,  $I_i(a) = H_i'(g_i^{-1}(a))$  and, in conjunction with (6), we have that  $I_i^{-1}(\delta) = h_i(\delta)/\delta$  as desired. ♠

With this lemma, the following "effective bandwidth" theorem for multiclass sources is immediate by independence.

*Theorem 1:* Assume that the arrival processes  $A_i$  all satisfy the conditions of the Gartner-Ellis Theorem and that the  $H_i$  are all strictly convex. For any  $\delta > 0$  and  $N \in \mathbf{C}$ ,

$$I(N, c) \geq \delta \Leftrightarrow \sum N_i \alpha_i(\delta) \leq c.$$

*Proof:* Let  $h$  be the log-moment generating function for the aggregate arrival process. Clearly,

$$h(\delta) = \sum N_i h_i(\delta).$$

Let the inverse of  $I(N, \cdot)$  be  $I_N^{-1}$ . Thus, by the argument in the preceding lemma,

$$I_N^{-1}(\delta) = \frac{h(\delta)}{\delta} = \frac{\sum N_i h_i(\delta)}{\delta} = \sum N_i \alpha_i(\delta)$$

as desired. ♠

This theorem shows that, under weak conditions on the arrival processes, effective bandwidths exist for the measure of congestion (3). The large deviations approach used is a unified framework to handle buffer sources modeled in different ways, as we shall see in the next section.

### III. MODELS OF ATM BUFFER SOURCES

We now consider several models of buffer sources used to characterize bursty ATM traffic. In each case, an expression for the effective bandwidth is found.

#### A. Constant Rate and Memoryless Sources

For sources with a constant arrival rate of  $R$  cells/second,  $A(t) = Rt$  for  $t > 0$ . Thus,  $h(\delta) = R\delta$ ,  $H(R) = 0$ , and  $H(M) = \infty$  for all  $M \neq R$ . Therefore, the hypothesis of Theorem 1 is satisfied and the effective bandwidth of this source is  $\alpha(\delta) = R$ . Note that, in the notation of Section II,  $\gamma = \Gamma = R$  for a constant rate source.

For memoryless (Poisson) sources with intensity  $R$  cells/second,  $h(t) = R(e^\delta - 1)$ . Thus,  $H(M) = M \log(M/R) - M + R$ . So, the hypothesis of Theorem 1 is satisfied and the effective bandwidth of this source is  $\alpha(\delta) = R(e^\delta - 1)/\delta$ . Note that  $\gamma = 0$  and  $\Gamma = \infty$  for a Poisson source.

#### B. Discrete-Time Markov Sources

We call a buffer source a discrete-time Markov source if there is a discrete-time Markov chain  $Z$  and a real constant  $R$  such that the number of arrivals to the buffer in interval of (continuous) time  $(nR^{-1}, (n+1)R^{-1})$  is a function of  $Z_n$ . We take the state space of  $Z$  to be  $1, 2, \dots, m$  and we let  $Q$  be its irreducible and aperiodic transition probability matrix. Let  $\Lambda_i$  be the number of cells that arrive in the interval  $(nR^{-1}, (n+1)R^{-1})$  when  $Z_n = i$ . We assume  $0 \leq \Lambda_i \leq \Lambda_{i+1} < \infty$  for all  $i = 1, \dots, m-1$ . Therefore, in the notation of Section II,  $\gamma = R\Lambda_1$ ,  $\Gamma = R\Lambda_m$ , and  $\bar{\gamma} := R \sum_i \pi_i \Lambda_i$  where  $\pi$  is the invariant of  $Q$ :  $\pi Q = \pi$ .

By an argument using the backward equation and Perron-Frobenius Theory [3],

$$h(\delta) = R \log [\rho(e^{\delta \Lambda} Q)] \quad (7)$$

where  $\Lambda = \text{diag}(\Lambda_1, \dots, \Lambda_m)$ , and  $\rho(F)$  is the spectral radius of matrix  $F$ .

$h$  is differentiable (and analytic) as a consequence of the perturbation theory of matrices (see [2, pp. 190-191]) and, therefore, satisfies the conditions of the Gartner-Ellis Theorem. In Section II, we established that  $h$  is convex. A simple consequence of [16, Lemma 3.4] is that either  $h$  is affine or strictly convex.

$h(0) = 0$  implies that the affine case is the constant rate source of Section III-A. If  $h$  is strictly convex, by direct calculation starting from (5) we get that  $H' = h'^{-1}$ . Thus,  $H'$  is strictly increasing which implies that  $H$  is strictly convex as well. So, the hypothesis of Theorem 1 is satisfied, and the effective bandwidth of this source is  $\alpha(\delta) = h(\delta)/\delta$ . This source is a special case of [3, Example 2.3], wherein the rates  $\Lambda_i$  are random.

*Two-State Discrete-Time Markov Source Example:* If the Markov chain considered is of the two-state ( $m = 2$ ) type, then by direct calculation,

$$h(\delta) = R \log \left[ \frac{1}{2} \left( a(\delta) + \sqrt{a^2(\delta) + 4b(\delta)} \right) \right]$$

where

$$a(\delta) = Q_{1,1} e^{\delta \Lambda_1} + Q_{2,2} e^{\delta \Lambda_2}$$

and

$$b(\delta) = e^{\delta(\Lambda_1 + \Lambda_2)} (1 - Q_{1,1} - Q_{2,2}).$$

#### C. Markov Fluids

A source is called a Markov fluid if its time derivative is a function of a continuous-time Markov chain on a finite-state space. As for the discrete-time Markov sources, we let  $1, \dots, m$  be the state space. Let  $Q$  be the irreducible transition rate matrix of the Markov fluid's time derivative and  $\Lambda_i$  be the arrival rate of cells when the time derivative of the Markov fluid is in state  $i$ . We make the same assumption on the parameters  $\Lambda_i$  that we made in the discrete-time Markov source case.

By an argument similar to that for discrete-time Markov sources (see the Appendix),

$$h(\delta) = \mu(Q + \delta \Lambda)$$

where  $\Lambda$  is defined and  $\mu(F)$  is the largest real eigenvalue of the matrix  $F$ . The same argument used for discrete-time Markov sources verifies that the hypothesis of Theorem 1 is satisfied.

*Two-State Markov Fluids Example:* If the Markov fluid considered is of the two-state ( $m = 2$ ) type, then by direct calculation

$$h(\delta) = \frac{1}{2} \left( -a(\delta) + \sqrt{a^2(\delta) - 4b(\delta)} \right)$$

where

$$a(\delta) = Q_{1,2} + Q_{2,1} - \delta(\Lambda_2 - \Lambda_1)$$

and

$$b(\delta) = \delta^2 \Lambda_2 \Lambda_1 - \delta(Q_{1,2} \Lambda_2 + Q_{2,1} \Lambda_1).$$

This is the effective bandwidth result in [12], [13].

#### D. Markov-Modulated Poisson Process

A source to a buffer is called a Markov-modulated Poisson process (MMPP) if the cell arrivals are Poisson with intensity  $\lambda$ , where  $\lambda$  is a function of a continuous-time Markov chain. We assume that the space  $\Lambda_1, \dots, \Lambda_m$  of intensities satisfies the conditions of the previous examples and that the transition rate matrix  $Q$  is irreducible.

By an argument similar to that for discrete-time Markov sources (again, see the Appendix)

$$h(\delta) = \mu(Q + (e^\delta - 1)\Lambda)$$

and the hypothesis of Theorem 1 is satisfied.

#### IV. CONCLUSIONS

Effective bandwidth results for the continuous-time Markovian sources of Sections III-C and III-D were also obtained in [8] using spectral decomposition methods [21], [9]. They found the same effective bandwidth formulas and establish (2) for buffers with multiclass Markov fluid sources and buffers with multiclass MMPP sources. The effective bandwidth results in Section II (using the large deviations approach) are more general than those of [8] and our measure of congestion (3) allows us to handle a buffer with sources modeled in different ways (e.g., a buffer with two sources: one modeled as a Markov fluid and the other as a MMPP). Recently, in [5], (2) was established for the stationary Lindley buffer process (discrete time) and they found an effective bandwidth result for a buffer using the simple "randomized priority" processor sharing rule [20].

In summary, we have shown the existence of effective bandwidths for a large class of sources commonly used to model ATM traffic. Given the effective bandwidths of a buffer's sources (i.e., the functions  $\alpha_i$  for the buffer of Section II), one can determine its spare capacity to accept more calls,  $c - \sum_{i=1}^K N_i \alpha_i(\delta)$ , which can be an integral part of network resource management [13], [14], [19].

#### APPENDIX

##### BACKWARD EQUATION APPROACH TO EVALUATE THE EFFECTIVE BANDWIDTH FOR MARKOVIAN SOURCES

For the Markov fluid source of Section III-C, let  $A(s, t)$  be the number of arrivals in the interval  $(s, t)$ ,  $x$  be the irreducible modulating Markov chain with rate matrix  $Q$  and invariant distribution  $\pi$ , and  $\psi_j(\delta, t) = E[\exp(\delta A(0, t)) | x(0) = j]$ .

The claim is that

$$h(\delta) := \lim_{t \rightarrow \infty} t^{-1} \log E \exp(\delta A(0, t)) = \mu(Q + \delta\Lambda).$$

To show this, we begin with a standard backward equation argument: for positive  $\epsilon \ll 1$ :

$$\psi_j(\delta, t) = E\left(E[e^{\delta A(0, t)} | x(\epsilon)] | x(0) = j\right) \quad (8)$$

$$= \sum_i \psi_i(\delta, t - \epsilon) e^{\epsilon Q(j, i)} e^{\epsilon \delta \Lambda_j} + o(\epsilon). \quad (9)$$

Since  $\exp(\epsilon Q)(j, i) = (I + \epsilon Q)(j, i) + o(\epsilon)$  and  $\exp(\epsilon \delta \Lambda_j) =$

$1 + \epsilon \delta \Lambda_j + o(\epsilon)$ , we get (after a little rearrangement)

$$\frac{\psi_j(\delta, t) - \psi_j(\delta, t - \epsilon)}{\epsilon} = \psi_j(\delta, t - \epsilon)(Q(j, j) + \delta \Lambda_j) + \sum_{i \neq j} \psi_i(\delta, t - \epsilon) Q(j, i) + O(\epsilon).$$

Letting  $\epsilon \rightarrow 0$ , we get

$$\frac{\partial}{\partial t} \psi_j(\delta, t) = \psi_j(\delta, t)(Q(j, j) + \delta \Lambda_j) + \sum_{i \neq j} \psi_i(\delta, t) Q(j, i).$$

In matrix form, this equation is

$$\frac{\partial}{\partial t} \Psi(\delta, t) = (Q + \delta \Lambda) \Psi(\delta, t)$$

where  $\Psi^T(\delta, t) = (\psi_1(\delta, t), \dots, \psi_m(\delta, t))$ . Thus,

$$\Psi(\delta, t) = \exp((Q + \delta \Lambda)t) \mathbf{1}$$

where  $\mathbf{1} = \Psi(\delta, 0)$  is a column of 1's.

Therefore,

$$h(\delta) = \lim_{t \rightarrow \infty} \frac{1}{t} \log (\pi^T \exp((Q + \delta \Lambda)t) \mathbf{1}).$$

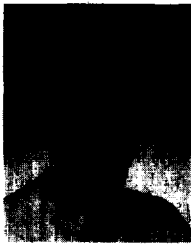
First note that  $\exp(Q + \delta \Lambda)$  is a nonnegative matrix (see [15, Exercise 6.5.4e and Theorems 6.2.9(g) and 6.2.38]). Choose  $a$  large enough such that  $aI + Q + \delta \Lambda \geq 0$ . This is possible since  $Q_{i,j} \geq 0$  for all  $i \neq j$ . Thus,  $\exp(Q + \delta \Lambda) = \exp(aI + Q + \delta \Lambda) \exp(-aI) \geq e^{-a} \exp(aI + Q + \delta \Lambda) \geq 0$ . Because of the irreducibility assumption, we can use the same Perron-Frobenius argument in [3] on the matrix  $\exp(Q + \delta \Lambda)$  to obtain  $h(\delta) = \log(\rho(\exp(Q + \delta \Lambda)))$ . The result then follows from  $\rho(\exp(F)) = \exp(\mu(F))$ , where  $\mu(F)$  is the largest eigenvalue of  $F$ .

For the case of the MMPP source of Section III-D, we use the fact that if  $\xi$  is a Poisson random variable with mean  $\epsilon \Lambda_j$ , then  $E \exp(\delta \xi) = \exp(\epsilon \Lambda_j (e^\delta - 1))$ . So, this argument will give us the formula for  $h$  in Section III-D by simply substituting the expression " $\exp(\epsilon \Lambda_j (e^\delta - 1))$ " for " $\exp(\epsilon \delta \Lambda_j)$ " in (9).

#### REFERENCES

- [1] D. Anick, D. Mitra, and M. M. Sondhi, "Stochastic theory of a data-handling system with multiple sources," *Bell Syst. Tech. J.*, vol. 61 no. 8, pp. 1871-1894, 1982.
- [2] J.A. Bucklew, *Large Deviation Techniques in Decision, Simulation and Estimation*. New York: Wiley, 1990.
- [3] C.-S. Chang, "Stability, queue length and delay. Part II: Stochastic queueing networks," in *Proc. IEEE CDC*, Tucson, AZ, 1992, pp. 1005-1010.
- [4] C. Courcoubetis and J. Walrand, "Note on effective bandwidth of ATM traffic," preprint.
- [5] G. de Veciana and J. Walrand, "Effective bandwidths: Call admission, traffic policing, and filtering for ATM networks," preprint, 1993.
- [6] R. Ellis, "Large deviations for a general class of random vectors," *Ann. Probab.*, vol. 12, pp. 1-12, 1984.
- [7] A.I. Elwalid and D. Mitra, "Fluid models for the analysis and design of statistical multiplexing with loss priorities on multiple classes of bursty traffic," in *Proc. IEEE INFOCOM* 1991, pp. 3c.4.1-3c.4.11.
- [8] A.I. Elwalid and D. Mitra, "Effective bandwidth of general Markovian traffic sources and admission control of high speed networks," in *Proc. IEEE INFOCOM*, 1993, vol. 1, pp. 3a.2.1-3a.2.10.
- [9] A.I. Elwalid, D. Mitra, and T.E. Stern, "Statistical multiplexing of Markov modulated sources: Theory and computational algorithms," in *Proc. ITC'13*, Copenhagen, Denmark, June 1991, pp. 495-500.

- [10] G. Gallassi, G. Rigolio, and L. Fratta, "ATM: Bandwidth assignment and bandwidth enforcement policies," in *Proc. GLOBECOM*, 1989, vol. 3, pp. 1788-1793.
- [11] J. Gartner, "On large deviations from invariant measure," *Theory Prob. Appl.*, vol. 22, pp. 24-39, 1977.
- [12] R.J. Gibbens and P.J. Hunt, "Effective bandwidths for multi-type UAS channel," *Queueing Syst.*, vol. 9, pp. 17-28, 1991.
- [13] R. Guerin, H. Ahmadi, and M. Naghshineh, "Equivalent capacity and its application to bandwidth allocation in high-speed networks," *IEEE J.Select. Areas Commun.*, vol. 9, no. 7, pp. 968-981, 1991.
- [14] R. Guerin and L. Gun, "A unified approach to bandwidth allocation and access control in fast packet-switched networks," in *Proc. IEEE INFOCOM*, 1992, vol. 1, pp. 1-12.
- [15] R.A. Horn and C.R. Johnson. *Topics in Matrix Analysis*. 1986.
- [16] I. Iscoe, P. Ney, and E. Nummelin, "Large deviations of uniformly recurrent Markov additive processes," *Adv. Appl. Prob.*, vol. 6, pp. 373-412, 1985.
- [17] F.P. Kelly, "Effective bandwidths of multi-class queues," *Queueing Syst.*, vol. 9, pp. 5-16, 1991.
- [18] G. Kesidis, "Cell loss estimation in high-speed digital networks," Ph.D. Dissert., EECS Dept., U.C. Berkeley, 1992.
- [19] G. Kesidis and J. Walrand, "Traffic policing and enforcement of effective bandwidth constraints in ATM networks," submitted to IEEE/ACM Trans. Netw.
- [20] J.M. Pitts and J.A. Schormans, "Analysis of ATM switch model with time priorities," *Electron. Lett.*, vol. 26, no. 15, pp. 1192-1193, July 1990.
- [21] T.E. Stern and A.I. Elwalid, "Analysis of separable Markov-modulated models for information handling systems," *Adv. Appl. Prob.*, Mar. 1991.



**George Kesidis** (S'91-M'92) was born in Toronto, Canada, in 1964. He received the B.A.Sc. degree in electrical engineering from the University of Waterloo in 1988, and the M.S. and Ph.D. degrees in electrical engineering from the University of California at Berkeley in 1990 and 1992, respectively.

He is currently an Assistant Professor in the Electrical and Computer Engineering Department of the University of Waterloo, Waterloo, Ont. His research interests include resource allocation, congestion control, and performance evaluation of high-speed networks.



**Jean Walrand** received the Ph.D. degree in electrical engineering from the University of California at Berkeley.

From 1979 to 1981, he taught at Cornell University and then joined the University of California at Berkeley, where he now is Professor of Electrical Engineering and Computer Science. His primary research interests are in queueing networks, stochastic control, and traffic control. He is the author of *An Introduction to Queueing Networks* (Prentice-Hall, 1988) and of *Communication Networks: A First Course* (Irwin/Aksen, 1991). He is Associate Editor for *Queueing Systems*, *Probability in the Engineering and Informational Sciences*, and *Mathematics of Operations Research*. He was awarded the Lanchester Prize by the Operations Research Society of America.



**Cheng-Shang Chang** (S'85-M'86-M'89-SM'93) received the B.S. degree in electrical engineering from the National Taiwan University, Taipei, Taiwan, in 1983, and the M.S. and Ph.D. degrees in electrical engineering from Columbia University, New York, NY, in 1986 and 1989, respectively.

Since 1989, he has been employed as a Research Staff Member at the IBM T.J. Watson Research Center, Yorktown Heights, NY. His current research interests are concerned with queueing theory, stochastic scheduling, and performance evaluation of telecommunication networks and parallel processing systems. He is an Associate Editor for *Operations Research*.