

Effective Bandwidth of General Markovian Traffic Sources and Admission Control of High Speed Networks

Anwar I. Elwalid, *Member, IEEE*, and Debasis Mitra, *Fellow, IEEE*

Abstract— The emerging high-speed networks, notably the ATM-based Broadband ISDN, are expected to integrate through statistical multiplexing large numbers of traffic sources having a broad range of burstiness characteristics. A prime instrument for controlling congestion in the network is admission control, which limits calls and guarantees a grade of service determined by delay and loss probability in the multiplexer. We show, for general Markovian traffic sources, that it is possible to assign a notional effective bandwidth to each source which is an explicitly identified, simply computed quantity with provably correct properties in the natural asymptotic regime of small loss probabilities. It is the maximal real eigenvalue of a matrix which is directly obtained from the source characteristics and the admission criterion, and for several sources it is simply additive. We consider both fluid and point process models and obtain parallel results. Numerical results show that the acceptance set for heterogeneous classes of sources is closely approximated and conservatively bounded by the set obtained from the effective bandwidth approximation. Also, the bandwidth-reducing properties of the Leaky Bucket regulator are exhibited numerically. For a source model of video teleconferencing due to Heyman *et al.* with a large number of states, the effective bandwidth is easily computed. The equivalent bandwidth is bounded by the peak and mean source rates, and is monotonic and concave with respect to a parameter of the admission criterion. Coupling of state transitions of two related asynchronous sources always increases their effective bandwidth.

I. INTRODUCTION

IN statistical multiplexing, which is the core of ATM-based IB-ISDN, we show that it is possible to assign a notional effective bandwidth to each source which reflects its characteristics, including burstiness, and the service requirements. The sources are given great generality; there are no restrictions on dimensions, homogeneity, or time reversibility. Yet, it is shown that the effective bandwidth is an explicitly identified quantity with provably correct asymptotic properties which can be obtained from simple and standard computations. In numerical evaluations of realistic admission control, approximations based on effective bandwidth perform very well. Importantly, the effective bandwidth of a source is independent of traffic submitted by other sources to the multiplexer. This fact makes the complexity of computing the effective bandwidth depend only on the source, not system, dimension; it also offers the promise of decentralized estimation from measurements and

enforcement of the effective bandwidth. Specifically, we show that the effective bandwidth of a Markovian source is the maximal real eigenvalue of a matrix, derived from the source parameters, network resources and service requirements, with dimension equal to the number of source states. Two parallel sets of results are obtained: one for a fluid model of statistical multiplexing with Markov-modulated fluid sources and the other for queues and point processes in which the traffic sources are, for example, Markov-modulated Poisson or phase renewal processes. The results extend the recent and important results of Gibbens and Hunt [20], who consider heterogeneous on/off fluid sources which alternate between exponentially distributed periods of transmission at the peak rate and quiescence. Even for the case of on/off fluid sources, the results here shed new light on the origins of key expressions in [20] and also on the expressions used there from earlier work by Anick *et al.* [2] and Kosten [25].

The imminence of new services with a broad range of burstiness characteristics and their integration through statistical multiplexing has focused attention on call admission as the prime instrument of rate-based congestion control. For a survey of issues, approaches, and analyses, see [33] and, for a more recent update, [34]. By preventing admission to an excessive number of calls or sources to the multiplexer, call admission policies strive to strike a balance between grade of service (as determined by delay and cell loss probability, for instance) and efficient use of network resources. Designs based on peak rates and mean rates are two extremes; hence, it is no accident that the effective bandwidth of a source is proven here to be bounded by these two rates. Designers have gravitated toward the concept of effective bandwidth because it promises simplicity and the hope that it might be a bridge to familiar circuit-switched network designs. It should be emphasized that the notion of effective bandwidth is intimately connected with admission control and the associated service requirements. Consequently, it is determined by the source characteristics in conjunction with the admission criteria. Gibbens and Hunt [20], Kelly [24], Guerin *et al.* [17], Kesidis and Walrand [27], Chang [7], and Whitt [40] offer different approaches to effective bandwidth. Kelly finds effective bandwidth for $GI/G/1$ queues. Guerin *et al.* independently obtain the formulas in [20] through insightful interpretations of the results in [2] and extend them through heuristics. Both [27] and [7] consider the general problem of the existence of an effective bandwidth of stationary and ergodic sources. Kesidis and Walrand take a

Manuscript received July 1992; revised February 1993; approved by IEEE/ACM TRANSACTIONS ON NETWORKING Editor Moshe Sidi.

The authors are with AT&T Bell Laboratories, Murray Hill, NJ 07974. (email: mitra@research.att.com)
IEEE Log Number 9211037.

large deviation approach to determine the effective bandwidth, where the admission criteria is identical to that in this paper. As is typical with this approach, they give the effective bandwidth in terms of the solution of a substantial variational problem; this problem is solved only for two-state on/off sources. There is an intimate connection between the behavior of tail probabilities of queue lengths and effective bandwidth. Whitt gives a detailed treatment of this connection for multiclass queues; see also Sohraby [38]. Chang develops bounding techniques for tail behavior of queues in networks and, among other results, explicitly connects the bounds for two-state on/off sources to the effective bandwidth in [20]. Other approaches to admission control and effective bandwidth based on loss networks are due to Hui [22] and Lindberger [28].

In the model of statistical multiplexing considered, each fluid source is characterized by $(\mathbf{M}, \boldsymbol{\lambda})$ where \mathbf{M} is the infinitesimal generator of a controlling Markov chain. The source generates fluid at the constant rate λ_s when in state s . The mean source rate is denoted by $\bar{\lambda}$ and the peak rate by $\hat{\lambda}$. The multiplexing buffer is serviced by a channel of constant capacity, or rate, c . Let $G(B)$ denote the stationary $\Pr[X \geq B]$ where X represents the random buffer content and interpret $G(B)$ to be the overflow probability for a buffer of size B . For given B and p , let the service requirement be $\{G(B) \leq p\}$, which is also taken to be the admission criterion. We think of p as being small, of the order of 10^{-9} .

Now consider the statistical multiplexing system in which there is only a single-source $(\mathbf{M}, \boldsymbol{\lambda})$. There are no restrictions on the dimensions of this source. We show that in the asymptotic regime where $p \rightarrow 0$ and $B \rightarrow \infty$ in such a manner that $\log p/B \rightarrow \zeta \in [-\infty, 0]$, the admission criterion is satisfied if $e < c$ and violated if $e > c$. We call e the effective bandwidth and show that it is the maximal real eigenvalue of the matrix $[\mathbf{A} - \frac{1}{\zeta} \mathbf{M}]$ where $\mathbf{A} = \text{diag}(\boldsymbol{\lambda})$. The effective bandwidth e depends on $(\mathbf{M}, \boldsymbol{\lambda})$, of course, and on the buffer and overflow probability only through ζ .

Now suppose that the single source just considered is, in fact, the aggregate of K arbitrary sources, $(\mathbf{M}^{(k)}, \boldsymbol{\lambda}^{(k)})$ ($1 \leq k \leq K$). We obtain a result of remarkable simplicity: the effective bandwidth $e = \sum e^{(k)}$, where $e^{(k)}$ is the effective bandwidth of the source $(\mathbf{M}^{(k)}, \boldsymbol{\lambda}^{(k)})$ computed as if it is a single source in the system.

In all important respects, the results carry over to the framework of queues and point processes. The source characterization differs only in that λ_s is the rate of the Poisson stream which is generated by the source when in state s . The effective bandwidth of the single-source $(\mathbf{M}, \boldsymbol{\lambda})$ in the multiplexing system is now the maximum real eigenvalue of $[\frac{1}{e\zeta} \mathbf{A} - \frac{1}{1-e\zeta} \mathbf{M}]$ for ζ defined as before. In the rest of the paper, we focus on the fluid model and handle the queueing model exclusively in Section VII.

We show that, in the fluid model, the effective bandwidth decreases monotonically with increasing ζ from $\hat{\lambda}$ at $\zeta = -\infty$ to $\bar{\lambda}$ at $\zeta = 0$. We also show that the coupling of state transitions from two asynchronous sources with identical infinitesimal generators and proportional rate vectors always leads to an increase in effective bandwidth. Examples show that this is not true, in general, if the rate vectors are arbitrary.

These facts are important if the pricing of network services is based on effective bandwidth.

The following is an observation on the effective bandwidth which may be useful for its estimation from measurements. Consider a testbed in which the source supplies a buffer which is serviced by a channel of variable capacity c . The effective bandwidth e is that value of c for which the asymptotic slope of $\log G(x)$ equals ζ .

The additive form in the effective bandwidth of K sources has simplifying consequences for the call admission problem with multiple heterogeneous classes of sources. We want $\mathcal{A}(B, p) = \{\mathbf{K} = (K^{(1)}, \dots, K^{(J)}) : G_{\mathbf{K}}(B) \leq p\}$. The asymptotic result is that $\mathcal{A}(B, p)$ is essentially the simplex $\sum e^{(j)} K^{(j)} < c$, where $e^{(j)}$ is the effective bandwidth of a single source of class j .

This asymptotic result motivates the approximation $\sum e^{(j)} K^{(j)} < c$ to the acceptance set in real, nonasymptotic cases. We have tested the goodness of this approximation for a variety of classes which display different burstiness aspects. Note that both the exactly calculated and the approximate acceptance sets are not exactly simplexes because of the integrality of $\{K^{(j)}\}$. It is our experience (reported in Section V) that the approximation is uniformly good provided that B is at least moderately large and, often, even when B is small. Also, we have observed that, importantly, the effective bandwidth approximation provides a conservative bound on the acceptance set.

Since so much emphasis has been given in prior work to two-state on/off sources, it is worth noting some of the reasons for considering higher-dimensional sources. First, as has been noted by Heyman *et al.* [21] for video teleconference traffic, two-state sources do not capture essential traffic features. In fact, Heyman *et al.* used models for individual sources which have over 600 states. See also Maglaris *et al.* [29] for other high-order video source models. Second, higher-order models are concomitant with the need to design call admission in conjunction with traffic monitoring and regulation. Elwalid and Mitra [11] have studied regulated traffic and its approximate Markovian characterization, which, in the case of the simplest class of regulators, has dimension one more than that of the original source model. Finally, even for on/off sources, there is considerable interest in the effects of variability of the on and off periods [5]. Such analysis requires higher-dimensional source models.

Results reported in Sections V and VI specifically address the points raised previously. In Section VI, we show that the effective bandwidth of a video teleconference source model with a large number of states, one recommended by Heyman *et al.*, can be quite easily calculated. Also, in Section V we numerically examine the bandwidth-reducing features of a Leaky Bucket regulator. These results are in agreement with Anantharaman and Konstantopoulos [1], from where it can be inferred that the effective bandwidth of the output of the Leaky Bucket is monotonic with respect to the regulator's parameter. Finally, Section V gives a simple source model with four states which accommodate hyperexponentially distributed on and off periods and also presents data on their influence on effective bandwidth.

The mathematical developments belong to two different categories: the first has to do with the analysis of just the single source, and the second with the algebraic decompositions which give the additive form to the effective bandwidth of several sources. The essential steps in the first category are broadening of the scope of the standard eigenvalue problem by introducing an inverse eigenvalue problem and investigating the growth properties of the maximum real eigenvalue with respect to a parameter in the problem. It is to the inverse problem that we bring to bear a fundamental result due to Cohen [8],[9] and Friedland [16] on the convex behavior of the maximum real eigenvalue of essentially nonnegative matrices with respect to all diagonal elements. In the second category, the algebraic theory which gives the important decompositions is based on Kronecker representations and separability, which has its antecedents in the work of Anick *et al.* [2], Kosten [25],[26], Mitra [30], Stern and Elwalid [35], and Elwalid *et al.* [15].

II. PRELIMINARIES

This section, which is in three parts, begins by giving some basic background facts about the statistical multiplexing system. Computation of the spectral expansion of the system's stationary distribution involves a standard eigenvalue problem. The second part of this section (Section II-B) points out that, in this paper, it will be necessary to broaden the scope of the eigenvalue problem by introducing a parameter (the channel capacity) and view the eigenvalues as functions of this parameter and, also, to look at the inverse problem, which turns out to be an eigenvalue problem as well. Finally, the last part of this section (Section II-C) presents some known facts about essentially nonnegative matrices and the maximal real eigenvalues critical for the analytic development in subsequent sections.

A. The Statistical Multiplexing System

The statistical multiplexing system consists of a buffer which is supplied by various statistically independent Markov-modulated fluid sources and serviced by a channel of constant capacity, i.e., rate c . For our purposes, it will suffice to lump the source description into a single aggregate Markov-modulated fluid source with state space \mathcal{S} and irreducible generator \mathbf{M} . (In Section IV, it will be necessary to consider the detailed structure of \mathbf{M} implied by the presence of several lower-order sources.) This aggregate source generates fluid at the constant rate λ_s when in state s ($s \in \mathcal{S}$). Let $\boldsymbol{\lambda} = \{\lambda_s | s \in \mathcal{S}\}$. Thus, the aggregate source is characterized by $(\mathbf{M}, \boldsymbol{\lambda})$. We also let the rate matrix $\mathbf{A} = \text{diag}(\boldsymbol{\lambda})$.

Let Σ and X denote the stationary aggregate-source state and buffer content, respectively. Let the stationary state distribution of the multiplexing system be denoted by $\boldsymbol{\pi}(x)$, where $\boldsymbol{\pi}(x) = \{\pi_s(x) | s \in \mathcal{S}\}$ and

$$\pi_s(x) \triangleq \Pr(\Sigma = s, X \leq x) \quad (s \in \mathcal{S}, 0 \leq x \leq \infty). \quad (1)$$

The governing system of differential equations is

$$\frac{d}{dx} \boldsymbol{\pi}(x) \mathbf{D} = \boldsymbol{\pi}(x) \mathbf{M} \quad (0 \leq x \leq \infty) \quad (2)$$

where $\mathbf{D} \triangleq \mathbf{A} - c\mathbf{I}$ and \mathbf{I} are the identity matrixes and the diagonal element $D_{ss} = (\lambda_s - c)$ is the drift, or rate of change, in the buffer content when the source is in state s . Hence, we call \mathbf{D} the drift matrix.

The stationary probability vector for the aggregate source is denoted by \mathbf{w} ; hence, $\mathbf{w}\mathbf{M} = \mathbf{0}$ and $\langle \mathbf{w}, \mathbf{1} \rangle = 1$. The symbol $\langle \cdot, \cdot \rangle$ denotes the inner product of vectors and $\mathbf{1}$ is the vector in which all elements are unity. The ergodicity condition is

$$\bar{\lambda} < c \quad (3)$$

where the mean source rate is

$$\bar{\lambda} \triangleq \langle \boldsymbol{\lambda}, \mathbf{w} \rangle. \quad (4)$$

We denote the peak source rate by $\hat{\lambda}$, i.e., $\hat{\lambda} \triangleq \max_s \lambda_s$. To rule out the trivial case in which there is never any accumulation in the multiplexing buffer, we assume that $c < \hat{\lambda}$.

Since the stationary state distribution is a bounded solution, it has the spectral representation

$$\boldsymbol{\pi}(x) = \sum_{i: \text{Re } z_i < 0} a_i \boldsymbol{\phi}_i e^{z_i x} + \mathbf{w} \quad (5)$$

where $(z_i, \boldsymbol{\phi}_i)$ is an eigenvalue/eigenvector pair. Such pairs are solutions to the eigenvalue problem

$$z \boldsymbol{\phi} \mathbf{D} = \boldsymbol{\phi} \mathbf{M}. \quad (6)$$

The eigenvalues with negative real parts are indexed as

$$0 > \text{Re } z_1 \geq \text{Re } z_2 \geq \text{Re } z_3 \geq \dots \quad (7)$$

If (as will turn out to be the case in all the cases considered in this paper) z_1 is real and $z_1 > \text{Re } z_i$ for all $i > 1$, then z_1 is called the *dominant eigenvalue*.

In the spectral expansion, the coefficients $\{a_i\}$ are obtained by solving a system of linear equations which are obtained by the following boundary conditions (see, for instance, [30])

$$D_{ss} > 0 \Rightarrow \pi(s, 0) = 0. \quad (8)$$

It is known that the number of such conditions exactly equals the number of eigenvalues with negative real parts.

Let the stationary buffer overflow distribution be given by $G(x)$, i.e.,

$$\begin{aligned} G(x) &= \Pr(X \geq x) \\ &= 1 - \langle \boldsymbol{\pi}(x), \mathbf{1} \rangle \\ &= \sum_{i \geq 1} a_i \langle \boldsymbol{\phi}_i, \mathbf{1} \rangle e^{z_i x}. \end{aligned} \quad (9)$$

If z_1 is the dominant eigenvalue, then

$$G(x) \sim a_1 \langle \boldsymbol{\phi}_1, \mathbf{1} \rangle e^{z_1 x} \quad \text{as } x \rightarrow \infty. \quad (10)$$

Note that

$$z_1 = \lim_{x \rightarrow \infty} \frac{\log G(x)}{x}. \quad (11)$$

Plots of $\log G(x)$ versus x approach linearity as x increases, and the slope approaches z_1 .

B. The Inverse Eigenvalue Problem

Consider the eigenvalue problem in (6)

$$z\phi(\mathbf{A} - c\mathbf{I}) = \phi\mathbf{M}. \quad (12)$$

It is convenient to extend the scope of the problem by considering c to be a variable parameter and the eigenvalues to be functions of c , $z(c)$. The inverse problem requires c to be obtained for given z . The key fact in this connection is that this inverse problem is also an eigenvalue problem. For, writing $c = g(z)$, we obtain from (12)

$$g(z)\phi = \phi\mathbf{A}(z), \quad (13)$$

where

$$\mathbf{A}(z) = \mathbf{A} - \frac{1}{z}\mathbf{M}. \quad (14)$$

That is, $g(z)$ is an eigenvalue of the matrix $\mathbf{A}(z)$ in which z is a parameter.

The inverse eigenvalue problem, its maximal real eigenvalue, and the behavior of this eigenvalue as a function of z will be important in the subsequent development.

C. Essentially Nonnegative Matrices

A real matrix with nonnegative elements off the main diagonal is called essentially nonnegative. The matrix $\mathbf{A}(z)$ in (14) is essentially nonnegative for real and negative z . Since \mathbf{M} is irreducible, so is $\mathbf{A}(z)$. By adding $\sigma\mathbf{I}$ to $\mathbf{A}(z)$, where

$$\sigma > \left[\max_i \left(\frac{1}{z}M_{ii} - \lambda_i \right) \right]^+, \quad (15)$$

we obtain a matrix which is nonnegative and eigenvalues which are the eigenvalues of $\mathbf{A}(z)$ shifted by σ . Thus, the Perron-Frobenius Theory [17],[10] applies to $[\mathbf{A}(z) + \sigma\mathbf{I}]$, and we can infer the following for matrix $\mathbf{A}(z)$ ($z < 0$).

FACT 1: There exists a real eigenvalue $g_1(z)$ of the matrix $\mathbf{A}(z)$ such that to $g_1(z)$ can be associated a real vector ϕ_1 , where $\phi_1 > 0$ (elementwise), and $\min_s \lambda_s < g_1(z) < \bar{\lambda}$. If $g(z)$ is any other eigenvalue then $\text{Re } g(z) < g_1(z)$. The eigenvalue $g_1(z)$ is simple. ■

The eigenvalue $g_1(z)$ is referred to as the *maximal real eigenvalue*. The maximal real eigenvalue of an essentially nonnegative matrix need not be the eigenvalue with the largest modulus. The upper and lower bounds on $g_1(z)$ in Fact 1 correspond to the maximum and minimum row sums of $\mathbf{A}(z)$.

A result due to Cohen [8, Theor. 1, Corol. 2] allows the lower bound in Fact 1 to be sharpened. Although Cohen's result is stated for nonnegative matrices, it is readily adapted to $\mathbf{A}(z)$ ($z < 0$). (Recall from (4) that $\bar{\lambda}$ is the mean source rate.)

FACT 2: $\bar{\lambda} \leq g_1(z)$. ■

We shall also need the following fundamental and important result due to Cohen [8],[9] and Friedland [16]:

FACT 3: The maximal real eigenvalue of $(\mathbf{A} + \Delta)$ is a strictly convex function of Δ , where \mathbf{A} is any irreducible, essentially nonnegative matrix and Δ is any diagonal matrix

with diagonal elements which are not all identical. That is,

$$\begin{aligned} & r((1-h)\mathbf{A} + h(\mathbf{A} + \Delta)) \\ & < (1-h)r(\mathbf{A}) + hr(\mathbf{A} + \Delta) \quad (0 < h < 1) \end{aligned} \quad (16)$$

where $r(\cdot)$ is the maximal real eigenvalue. This result is equivalent to the positive definiteness of the Hessian $\mathbf{H} = \{H_{ij}\}$, where $H_{ij} = \partial^2 r(\mathbf{A} + \Delta) / \partial \Delta_{ii} \partial \Delta_{jj}$. ■

The antecedents of this result are rich and varied. Cohen [8] showed weak convexity by a Feynman-Kac formula for the maximal real eigenvalue of nonnegative matrices. Friedland [16], who was the first to show strict convexity, used a variational characterization by Donsker and Varadhan for the maximum real eigenvalue. Cohen [9] also showed strict convexity by using the Trotter product formula and a theorem on log convexity due to Kingman [23] as extended by Seneta [36].

Finally, we shall need recourse to a well-known result on nonnegative matrices [17],[10].

FACT 4: The maximal real eigenvalue of a nonnegative, irreducible matrix increases when any matrix element increases. ■

This result remains intact when the nonnegativity of the matrix is substituted by essential nonnegativity.

Note that increasing z in $\mathbf{A}(z) = [\mathbf{A} - \frac{1}{z}\mathbf{M}]$ has the effect of increasing the nonzero off-diagonal elements, and of decreasing the diagonal elements of the matrix. Hence, this result by no means implies the monotonicity of the maximal real eigenvalue $g_1(z)$ with respect to z . This topic, which is of central importance in our study, is examined in the next section.

III. A SINGLE SOURCE: MONOTONICITY OF EIGENVALUES AND EFFECTIVE BANDWIDTH

This section on the single source plays a pivotal role in this discussion. First, when the single source in the multiplexing system is allowed to be of arbitrary dimension (as it is in this section), then it can be construed to be the aggregate of many lower-order sources and the many-source problem becomes an extension of the single-source problem in which the new element is the algebraic exploitation of the structure implied by the presence of many sources. Such an extension is undertaken in the next section. Second, the qualitative properties of the eigenvalues, such as monotonicity and convexity, are established in Section III-A. Third, the asymptotic view of the admission control problem is introduced in Section III-B. The result identifying the effective bandwidth of the source as the maximal real eigenvalue of a simple matrix is proven there. Finally, in Section III-C we show that the effective bandwidth is a monotonic increasing and convex function of all state-dependent rates of the source. A corollary to this result is that, whenever we couple the state transitions of two sources having identical generators for their controlling Markov chains and proportional rate vectors, the effect is to increase the effective bandwidth.

Let the source be characterized by (\mathbf{M}, λ) where \mathbf{M} is any irreducible infinitesimal generator. The number of states in the controlling Markov chain, which is also the dimension

of \mathbf{M} and λ , is N . The system considered in this section consists of this source supplying a buffer which is serviced by a channel of capacity (rate) c . The admission control problem is to characterize sources for which the admission criterion $\{G(B) \leq p\}$ is satisfied.

A. Monotonicity of the Maximal Real Eigenvalue of the Inverse Problem

We examine the maximal real eigenvalue of the inverse eigenvalue problem (13). Recall that, in this problem, the parameter is z and the eigenvalue is $g(z)$:

$$g(z)\phi = \phi\mathbf{A}(z) \quad (17)$$

where

$$\mathbf{A}(z) = \mathbf{A} - \frac{1}{z}\mathbf{M}. \quad (18)$$

Making use of Fact 1 of Section II-C, the solutions to (17) for $z < 0$ are indexed thus:

$$g_1(z) > \text{Re } g_2(z) \geq \text{Re } g_3(z) \geq \dots \quad (19)$$

The maximal real eigenvalue is $g_1(z)$. From Facts 1 and 2,

$$\bar{\lambda} \leq g_1(z) \leq \hat{\lambda} \quad (z < 0). \quad (20)$$

We will find it convenient more than once to complement (17) by the form which is obtained by multiplying (17) by $(-z)$:

$$\{-zg(z)\}\phi = \phi[(-z)\mathbf{A} + \mathbf{M}]. \quad (21)$$

The matrix $(-z)\mathbf{A}(z)$ on the right side remains essentially nonnegative and hence has a maximal real eigenvalue, which we denote by $r(z)$. Note that $r(z) = -zg_1(z)$.

The form in (21) is useful for obtaining the limiting value of $g_1(z)$ as $z \rightarrow 0$. From Fact 1 of Section II-C, we know that $g_1(z)$ and $r(z)$ are simple; consequently, standard perturbation analysis applies [39]. Expanding $r(z)$ and ϕ in power series, $r(z) = r_0 + r_1z + r_2z^2 + \dots$, and $\phi(z) = \phi_0 + z\phi_1 + z^2\phi_2 + \dots$, we obtain from (21) $r_0 = 0$, $\phi_0 = \mathbf{w}$, and $r_1 = -\langle \lambda, \mathbf{w} \rangle = -\bar{\lambda}$. Hence,

$$g_1(0) = -\left[\lim_{z \rightarrow 0} \frac{r(z)}{z}\right] = -r_1 = \bar{\lambda}. \quad (22)$$

When $z \rightarrow -\infty$, it is apparent from (17) and (18) that $g_1(z) \rightarrow \hat{\lambda}$. To recapitulate, we have that the maximal real eigenvalue $g_1(z)$ of $\mathbf{A}(z)$ satisfies, for $z < 0$,

Proposition III.1

$$g_1(0) = \bar{\lambda} \quad \text{and} \quad g_1(-\infty) = \hat{\lambda}. \quad (23)$$

Hence, the bounds in (20) are tight.

The effect of decreasing z is to increase all the diagonal elements of $[(-z)\mathbf{A} + \mathbf{M}]$ for which the corresponding diagonal elements of \mathbf{A} are nonzero, while all other diagonal elements and all off-diagonal elements are not affected. Hence, it follows from Fact 4 of Section II-C that

$$\frac{\partial r}{\partial z} < 0 \quad (z < 0). \quad (24)$$

Next, we use Fact 3 to establish the convexity of $r(z)$. Let $z_2 < z_1 < 0$ and $0 < h < 1$. In Fact 3, identify \mathbf{A} with $[(-z_1)\mathbf{A} + \mathbf{M}]$ and Δ with $(z_1 - z_2)\mathbf{A}$ to obtain

$$r\{(1-h)z_1 + hz_2\} < (1-h)r(z_1) + hr(z_2). \quad (25)$$

This condition is equivalent to

$$\frac{\partial^2 r}{\partial z^2} > 0 \quad (z < 0). \quad (26)$$

In summary, we have from (24) and (26)

Proposition III.2 The maximal real eigenvalue $r(z)$ of the essentially nonnegative matrix $(-z)\mathbf{A}(z) = [(-z)\mathbf{A} + \mathbf{M}]$ ($z < 0$) is a monotonically decreasing convex function. Moreover, $r(z) \sim \hat{\lambda}|z|$ as $z \rightarrow -\infty$, and $r(0) = 0$. ■

Recalling that $r(z) = -zg_1(z)$, it follows from (24) and (26) that, for $z < 0$,

$$g_1(z) + zg_1'(z) > 0, \quad (27)$$

$$2g_1'(z) + zg_1''(z) < 0. \quad (28)$$

We are ready to prove one of our main results.

Proposition III.3 The maximal real eigenvalue $g_1(z)$ of the essentially nonnegative matrix $\mathbf{A}(z) = [\mathbf{A} - \frac{1}{z}\mathbf{M}]$ is monotonic, decreasing with increasing z ,

$$g_1'(z) < 0 \quad (z < 0). \quad (29)$$

Proof: From (28), $g_1'(z) < 0$ when $|z|$ is small, and from (27), $g_1'(z) < 0$ when $|z|$ is large. Now suppose that there exists some z for which (29) is false. Then, there exists intervals in which the sign of $g_1'(z)$ is uniform, and in neighboring intervals the signs are opposite. In particular, there must exist a common endpoint to two such contiguous intervals, say z_1 ($-\infty < z_1 < 0$), where a local maximum is reached, i.e.,

$$g_1'(z_1) = 0, \text{ and } g_1''(z_1) \leq 0. \quad (30)$$

Notice that $2g_1'(z_1) + z_1g_1''(z_1) \geq 0$, which contradicts (28). ■

We can also show that $g_1(z)$ is a concave function; the proof is omitted.

Standard perturbation analysis readily yields an expression for $g_1'(z)$. In (22), let z be perturbed to $z + \epsilon$ and consider an expansion in powers of ϵ , $g_1(z + \epsilon) = g_1(z) + \sum \epsilon^i g_{1i}(z)$, and similar expansions for the left and right eigenvectors. By equating coefficients of ϵ^0 and ϵ^1 , we obtain

$$g_1'(z) = \frac{1}{z^2} \frac{\phi(z)\mathbf{M}\psi(z)}{\phi(z)\psi(z)} \quad (31)$$

where $\phi(z)$ and $\psi(z)$ are, respectively, the left (row) and right (column) real eigenvectors of $\mathbf{A}(z)$ corresponding to the eigenvalue $g_1(z)$.

Notice that when the Markov chain with generator \mathbf{M} is time reversible, then \mathbf{M} is essentially symmetric and negative semi-definite. Also, the form in (31) immediately shows that $g_1'(z) < 0$ for $z < 0$ ($g_1'(0)$ is more delicate). This fact was established by Stern and Elwalid [35].

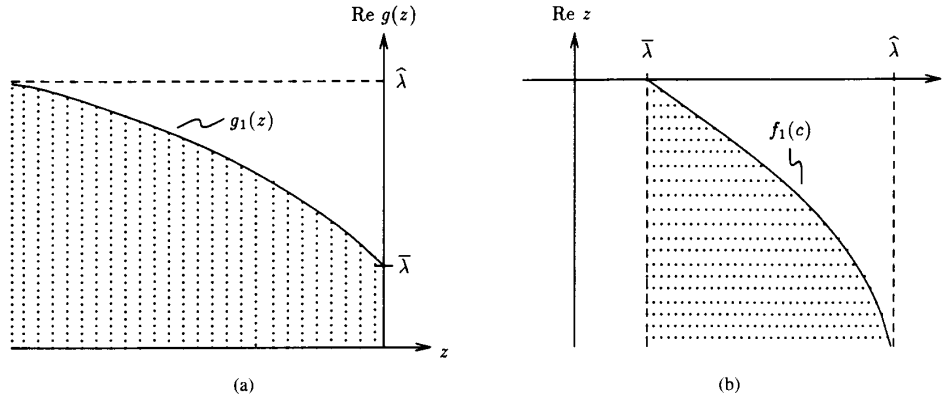


Fig. 1. (a) The shaded region contains the real parts of all solutions $g(z)$ to the inverse eigenvalue problem (17). (b) The shaded region contains the real parts of all solutions z with negative real parts to the eigenvalue problem (6).

An on/off source with exponentially distributed on and off periods is obtained by setting $\lambda_1 = 0$ in the following representation:

$$\mathbf{M} = \begin{bmatrix} -\alpha & \alpha \\ \beta & -\beta \end{bmatrix} \quad \text{and} \quad \boldsymbol{\lambda} = [\lambda_1 \quad \lambda_2]. \quad (32)$$

The reader may verify that

$$g_1(z) = \frac{1}{2z} [(\lambda_1 + \lambda_2)z + \beta + \alpha] - \frac{1}{2z} \sqrt{\{(\lambda_1 + \lambda_2)z + \alpha + \beta\}^2 - 4(\lambda_1 \lambda_2 z^2 + \beta \lambda_1 z + \alpha \lambda_2 z)} \quad (33)$$

This expression is central to Anick *et al.* [2] and Kosten [25]. Direct differentiation of (33), the procedure followed by Gibbens and Hunt [20], confirms that $g_1'(z) < 0$ for $z < 0$.

Fig. 1(a) is a sketch incorporating the results in Propositions III.1 and III.3.

The key duality between the two extremal eigenvalues of interest in the direct and inverse eigenvalue problems is now given.

Proposition III.4 For $c \in (\bar{\lambda}, \hat{\lambda})$ the dominant eigenvalue z_1 is the unique solution in $(-\infty, 0)$ satisfying

$$g_1(z_1) = c, \quad (34)$$

i.e., the dominant eigenvalue is the unique parameter in $[A - \frac{1}{z}\mathbf{M}]$ for which the maximal real eigenvalue is c .

Proof: Since $g_1(z)$ is monotonic, strictly decreasing for $z \in [-\infty, 0]$, and takes values between $[\bar{\lambda}, \hat{\lambda}]$, (34) has a unique solution. If z_2 is any other real solution to $g(z) = c$, then $z_2 < z_1$. If $z_2 > z_1$, then

$$c = g(z_2) \leq g_1(z_2) < g_1(z_1) = c$$

a contradiction. It only remains to show that the dominant eigenvalue is real. The proof is similar to that used to show that the maximal real eigenvalue in the inverse eigenvalue problem is real, and we omit the detailed proof. ■

Denote the unique inverse of g_1 in $[-\infty, 0]$ by f_1 , i.e.,

$$f_1(g_1(z)) = z \quad (z < 0).$$

Hence, f_1 maps $[\bar{\lambda}, \hat{\lambda}]$ to $[0, -\infty]$ and

$$z_1 = f_1(c) \quad (\bar{\lambda} \leq c \leq \hat{\lambda}). \quad (35)$$

It is easily seen from an application of the chain rule that

Proposition III.5

$$\frac{df_1}{dc} < 0 \quad (\bar{\lambda} < c < \hat{\lambda}) \quad (36)$$

i.e., the dominant eigenvalue z_1 is a monotonic, strictly decreasing function of the channel capacity c for $c \in (\bar{\lambda}, \hat{\lambda})$. ■

These results are incorporated in Fig. 1(b).

B. Asymptotics: Small Overflow Probabilities, Large Buffers

We now consider the admission control problem for an asymptotic regime in which the buffer overflow probability (p) is small, say of the order of 10^{-9} . For the scaling in this asymptotic regime to be meaningful, it is of course necessary to have large buffers. Enough is already known [see (10)] about the qualitative manner in which overflow probabilities scale with buffer size (B) to arrive at the following natural asymptotic regime, which is also the regime considered by Gibbens and Hunt [20]. Let $B \rightarrow \infty$, and also $p \rightarrow 0$, in such a manner that

$$\log p = \zeta B + O(1), \quad (37)$$

where $\zeta \in [-\infty, 0]$ is any $O(1)$ parameter. Hence, $\log p/B \rightarrow \zeta$.

Since in this section we are considering a system with a single source, the multiplexing is nonexistent. The problem here is to characterize sources which supply a system with a buffer of size B and channel capacity c , and for which the buffer overflow probability $G(B)$ does not exceed p in the aforementioned asymptotic regime.

Proposition III.6 Let the admission criterion be $G(B) \leq p$. Suppose $B \rightarrow \infty$ and $p \rightarrow 0$ in such a manner that $(\log p/B) \rightarrow \zeta \in [-\infty, 0]$.

If $g_1(\zeta) < c$, then the admission criterion is satisfied.

If $g_1(\zeta) > c$, then the admission criterion is violated, where $g_1(\zeta)$ is the maximal real eigenvalue of $\mathbf{A}(\zeta) = \mathbf{A} - \frac{1}{\zeta}\mathbf{M}$.

Proof: From (9),

$$G(B) = \sum_{i \geq 1} a_i \langle \phi_i, \mathbf{1} \rangle e^{z_i B}. \quad (38)$$

Recall that z_1 is the dominant eigenvalue, so that $z_1 > \operatorname{Re} z_i$ for all $i > 1$ and [see (34)] $c = g_1(z_1)$. So,

$$\frac{G(B)}{p} = a_1 \langle \phi_1, \mathbf{1} \rangle e^{(z_1 - \zeta)B} [1 + o(1)] \quad \text{as } pB \rightarrow \infty. \quad (39)$$

Now, from Proposition III-C, $g_1(z)$ decreases as z increases; hence, if $g_1(\zeta) < c$, then $z_1 < \zeta$ [see Fig. 1(a)] and $\{G(B)/p\} \rightarrow 0$ as $B \rightarrow \infty$. Therefore, the admission criterion is satisfied. Similarly, if $g_1(\zeta) > c$ then $z_1 > \zeta$ and $\{G(B)/p\} \rightarrow \infty$, so the admission criterion is violated. ■

This result justifies the use of the term “effective bandwidth” for the quantity $g_1(\zeta)$. We let $e = e(\mathbf{M}, \boldsymbol{\lambda}; B, p)$ denote the effective bandwidth of the source $(\mathbf{M}, \boldsymbol{\lambda})$ in the system for which the admission criterion is $G(B) \leq p$. That is,

$$e(\mathbf{M}, \boldsymbol{\lambda}; B, p) = g_1(\zeta) \quad (40)$$

where $g_1(\zeta)$ is the maximal real eigenvalue of the matrix $(\mathbf{A} - \frac{1}{\zeta} \mathbf{M})$ and $\zeta = \log p/B$.

The fact that B and p determine e only through ζ is a consequential fact which simplifies and benefits the design process (see the discussion in [20]). Note that the effective bandwidth is independent of the channel capacity. From (23), (29), and (40), we observe that the effective bandwidth decreases monotonically with increasing ζ from the peak source rate $\hat{\lambda}$ when $\zeta = -\infty$ to the mean source rate $\bar{\lambda}$ when $\zeta = 0$.

There are several effective numerical algorithms for calculating the maximal real eigenvalue and the Perron root of nonnegative matrices, such as the inverse iteration method (see [39] and [32]). Recall from (15), $(\mathbf{A} - \frac{1}{\zeta} \mathbf{M} + \sigma \mathbf{I})$ is a nonnegative matrix for $\sigma > \left[\max_i \left(\frac{1}{\zeta} \mathbf{M}_{ii} - \lambda_i \right) \right]^+$.

The discussion in Section II-B on the inverse eigenvalue problem indicates an interpretation of $g_1(\zeta)$ which may be quite useful for obtaining the effective bandwidth from measurements. Consider a testbed in which the source supplies a buffer which is emptied by a channel of (variable) capacity c . The effective bandwidth e is that value of c for which the asymptotic slope of $\log G(x)$ equals ζ .

C. Monotonicity and Convexity of the Effective Bandwidth with Respect to Source Rates

Here we investigate the influence of the source rates $\boldsymbol{\lambda} = (\lambda_1, \lambda_2, \dots, \lambda_N)$ on the effective bandwidth of the source, $e(\mathbf{M}, \boldsymbol{\lambda}; B, p)$. First, we establish that with \mathbf{M} , B , and p held fixed, the effective bandwidth is strictly monotonic (increasing with each increasing λ_i) and also convex in $\lambda_1, \lambda_2, \dots, \lambda_N$. Next, these properties are used to obtain the following inequality:

$$e(\mathbf{M}, \boldsymbol{\lambda}; B, p) > e(\mathbf{M}, a\boldsymbol{\lambda}; B, p) + e(\mathbf{M}, (1-a)\boldsymbol{\lambda}; B, p) \quad (41)$$

for all $a \in (0, 1)$. The right quantity is the sum of the effective bandwidths of two sources which have identical controlling generator \mathbf{M} for their Markov chains, proportional rate vectors,

and are statistically independent, i.e., asynchronous. The left quantity is the effective bandwidth of the source obtained by coalescing the two sources by coupling state transitions. The next section will show that the combined effective bandwidth of the two asynchronous sources is simply the sum of their individual effective bandwidths. Hence the result in (41) shows that coupling always increases the effective bandwidth. This result is important for admission control, traffic shaping, and policing, and in the use of the effective bandwidth concept for determining prices of network services.

Since in this subsection \mathbf{M} , B , and p are held fixed and only $\boldsymbol{\lambda}$ is varied, it is convenient to write for the effective bandwidth,

$$e(\boldsymbol{\lambda}) = g_1(\boldsymbol{\lambda}). \quad (42)$$

It is understood that in this subsection $g_1(\boldsymbol{\lambda})$ denotes the maximal real eigenvalue of $\left[\mathbf{A} - \frac{1}{\zeta} \mathbf{M} \right]$, where $\zeta = \log p/B$. Note that $e(0) = g_1(0) = 0$.

Proposition III.7 $e(\boldsymbol{\lambda})$ is monotonic, increasing when any rate λ_i increases. Also, $e(\boldsymbol{\lambda})$ is convex in $\lambda_1, \lambda_2, \dots, \lambda_N$.

Proof: Strict monotonicity follows from Fact 4 of Section II-C since $\left[\mathbf{A} - \frac{1}{\zeta} \mathbf{M} \right]$ is an essentially nonnegative and irreducible matrix, and increasing any λ_i increases a diagonal element of the matrix. Strict convexity follows from Fact 3 and a similar observation. ■

An immediate implication of the convexity of $e(\boldsymbol{\lambda})$ is that, for any two nonnegative (elementwise) and nonnull rate vectors $\boldsymbol{\lambda}_1$ and $\boldsymbol{\lambda}_2$,

$$e(a\boldsymbol{\lambda}_1 + (1-a)\boldsymbol{\lambda}_2) \leq ae(\boldsymbol{\lambda}_1) + (1-a)e(\boldsymbol{\lambda}_2) \quad (0 < a < 1) \quad (43)$$

with equality holding if and only if all elements of $(\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}_2)$ are identical.

Yet another implication is

Proposition III.8 For all $a \in (0, 1)$,

$$e(\boldsymbol{\lambda}) > e(a\boldsymbol{\lambda}) + e((1-a)\boldsymbol{\lambda}). \quad (44)$$

Proof:

$$\begin{aligned} e(\boldsymbol{\lambda}) - e(a\boldsymbol{\lambda}) &= \int_0^1 g_1'(a\boldsymbol{\lambda} + \tau(1-a)\boldsymbol{\lambda}) d\tau \\ &> \int_0^1 g_1'(\tau(1-a)\boldsymbol{\lambda}) d\tau \\ &= g_1((1-a)\boldsymbol{\lambda}) \\ &= e((1-a)\boldsymbol{\lambda}). \quad \blacksquare \end{aligned}$$

Note that arbitrary splittings of the source with rate vector $\boldsymbol{\lambda}$ into two asynchronous sources with rates $\boldsymbol{\lambda}_1$ and $\boldsymbol{\lambda}_2$ ($\boldsymbol{\lambda}_1 + \boldsymbol{\lambda}_2 = \boldsymbol{\lambda}$) do not generally preserve the inequality in (44). In fact, the following example shows that the reverse is not uncommon. Consider on/off sources with exponentially distributed on and off periods, for which the generator \mathbf{M} is given in (32). Let the source rate vector $\boldsymbol{\lambda} = (r \quad r)$, in which case the effective bandwidth calculated from (33) is r . Now consider

$$\boldsymbol{\lambda}_1 = (0 \quad r) \quad \text{and} \quad \boldsymbol{\lambda}_2 = (r \quad 0). \quad (45)$$

The reader may verify that

$$e(\boldsymbol{\lambda}_1 + \boldsymbol{\lambda}_2) < e(\boldsymbol{\lambda}_1) + e(\boldsymbol{\lambda}_2) \quad (46)$$

for all $\zeta < -|\beta - \alpha|/r$.

IV. MULTIPLE SOURCES

We extend the results in the preceding section to multiplexing systems with several sources. First, we consider K arbitrary Markov-modulated fluid sources and seek a characterization of the sources for which the admission criterion $\{G(B) \leq p\}$ is satisfied. As in the last section, the framework is asymptotic and the natural scaling in (37) is used. A key element of the asymptotic analysis here, as in Section III-B, is the monotonicity of the maximal real eigenvalue with respect to the parameter in the inverse eigenvalue problem. The essential new element here is the simple additive form of the equation having K terms (called here the ‘‘coupled eigenvalue problem’’), which is satisfied by the eigenvalues of the system. This, together with the accompanying result which represents the system eigenvector as the Kronecker product of K low-order eigenvectors, constitutes a major decomposition of the eigenvalue problem. The algebraic theory which gives the decomposition is based on Kronecker representations and separability, and its antecedents are the results in Anick *et al.* [2], Kosten [25],[26], Mitra [30] and Stern and Elwalid [35].

A. Representations

We suppose that there are K sources characterized by $(\mathbf{M}^{(k)}, \lambda^{(k)})$ ($k = 1, 2, \dots, K$). Assume that for every k , source k has $N^{(k)}$ states and the generator $\mathbf{M}^{(k)}$ is irreducible. Let $\mathbf{A}^{(k)} = \text{diag}(\lambda^{(k)})$. $\mathcal{S}^{(k)} = \{1, 2, \dots, N^{(k)}\}$ is the state space of source k .

The aggregate source is a continuous-time Markov chain with state space $\mathcal{S} = \{\mathbf{s} | \mathbf{s} = (s^{(1)}, \dots, s^{(K)}), s^{(k)} \in \mathcal{S}^{(k)}, 1 \leq k \leq K\}$. The states of the sources are statistically independent, and consequently the infinitesimal generator of the aggregate source is \mathbf{M} , where

$$\mathbf{M} = \mathbf{M}^{(1)} \otimes \mathbf{I} \otimes \dots \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{M}^{(2)} \otimes \mathbf{I} \otimes \dots \otimes \mathbf{I} + \dots + \mathbf{I} \otimes \dots \otimes \mathbf{I} \otimes \mathbf{M}^{(K)} \quad (47)$$

and \otimes denotes the Kronecker product. The Appendix gives information on definitions and results of Kronecker algebra which are used in this paper. The standard compact representation of the form in (47) is

$$\mathbf{M} = \mathbf{M}^{(1)} \oplus \mathbf{M}^{(2)} \oplus \dots \oplus \mathbf{M}^{(K)}, \quad (48)$$

a K -fold Kronecker sum. The generator \mathbf{M} is also irreducible.

The stationary probability vector of the aggregate source \mathbf{w} is the Kronecker product of the stationary probability vectors of the individual sources. That is, $\mathbf{w}\mathbf{M} = \mathbf{0}$ and $\langle \mathbf{w}, \mathbf{1} \rangle = 1$ where

$$\mathbf{w} = \mathbf{w}^{(1)} \otimes \mathbf{w}^{(2)} \otimes \dots \otimes \mathbf{w}^{(K)} \quad (49)$$

and $\mathbf{w}^{(k)}\mathbf{M}^{(k)} = \mathbf{0}$, $\langle \mathbf{w}^{(k)}, \mathbf{1} \rangle = 1$. The system rate matrix \mathbf{A} is

$$\mathbf{A} \triangleq \mathbf{A}^{(1)} \oplus \mathbf{A}^{(2)} \oplus \dots \oplus \mathbf{A}^{(K)}. \quad (50)$$

The system drift matrix \mathbf{D} is

$$\mathbf{D} \triangleq \mathbf{A} - c\mathbf{I}. \quad (51)$$

The ergodicity condition is $\bar{\lambda} < c$, where the mean rate of the aggregate source

$$\bar{\lambda} \triangleq \sum_k \bar{\lambda}^{(k)} = \sum_k \langle \lambda^{(k)}, \mathbf{w}^{(k)} \rangle. \quad (52)$$

The peak rate of the aggregate source

$$\hat{\lambda} \triangleq \sum_k \hat{\lambda}^{(k)}. \quad (53)$$

To avoid trivialities, we assume that $c < \hat{\lambda}$.

B. Decompositions

Let us transform the eigenvalue equation $z\phi\mathbf{D} = \phi\mathbf{M}$ to obtain the form of the inverse eigenvalue equation:

$$c\phi = \phi\mathbf{A}(z) \quad (54)$$

where

$$\mathbf{A}(z) \triangleq \mathbf{A} - \frac{1}{z}\mathbf{M}. \quad (55)$$

A key observation is that $\mathbf{A}(z)$, like \mathbf{A} and \mathbf{M} , also has the Kronecker sum form, which the reader is invited to verify:

$$\mathbf{A}(z) = \mathbf{A}^{(1)}(z) \oplus \mathbf{A}^{(2)}(z) \oplus \dots \oplus \mathbf{A}^{(K)}(z) \quad (56)$$

where,

$$\mathbf{A}^{(k)}(z) \triangleq \mathbf{A}^{(k)} - \frac{1}{z}\mathbf{M}^{(k)} \quad (1 \leq k \leq K). \quad (57)$$

From the eigenvalue and eigenvector results of Kronecker sums which are stated in the Appendix, we obtain

Proposition IV.1 A necessary and sufficient condition for (c, ϕ) to be a solution to the eigenvalue problem in (54) and thus also for (z, ϕ) to be a solution to the eigenvalue problem in (6) is

$$\left. \begin{aligned} g^{(k)}(z)\phi^{(k)}(z) &= \phi^{(k)}(z)\mathbf{A}^{(k)}(z) \quad (1 \leq k \leq K) \\ \sum_k g^{(k)}(z) &= c \end{aligned} \right\} \quad (58a) \quad (58b)$$

where the eigenvector

$$\phi = \phi^{(1)} \otimes \phi^{(2)} \otimes \dots \otimes \phi^{(K)}. \quad (59)$$

We have called (58) the ‘‘coupled eigenvalue problem’’ since it is a system of K eigenvalue problems in which the dimensions are only $N^{(k)}$ ($1 \leq k \leq K$), and (58b) couples the constituent problems. As an alternative to this formal approach, the reader is invited to postulate the form for the eigenvector in (59) and to verify that the eigenvalue equation is satisfied if (58) holds. The remaining necessity part of the proof consists of verifying that the right number of eigenvalues are obtained by this procedure. For given z and k , there are $N^{(k)}$ solutions to (58a). Denote these by $g_{i^{(k)}}^{(k)}(z)$, where $i^{(k)} \in \{1, 2, \dots, N^{(k)}\}$. Hence, (58b) is equivalent to the family of equations

$$\sum_{k=1}^K g_{i^{(k)}}^{(k)}(z) = c \quad (60)$$

in which all combinations of the subscripts are to be considered.

C. Dominant Eigenvalue, Maximal Real Eigenvalue

On examining the individual equations in (58a), we see that they are in the form of inverse eigenvalue problems, which are the subject of detailed investigation in Section III-A. It is known [see (19)] that, for $z < 0$, there exists a simple real solution $g_1^{(k)}(z)$, called the maximal real eigenvalue, such that

$$g_1^{(k)}(z) > \text{Re } g_2^{(k)}(z) \geq \text{Re } g_3^{(k)}(z) \geq \dots \quad (z < 0). \quad (61)$$

Moreover, it has been established in Propositions 22 and 28 that $g_1^{(k)}(z)$ monotonically decreases from $\widehat{\lambda}^{(k)}$ to $\overline{\lambda}^{(k)}$ as z increases from $-\infty$ to 0.

The analog of Proposition III.4 is

Proposition IV.2 For $c \in (\overline{\lambda}, \widehat{\lambda})$, the dominant eigenvalue z_1 is the unique solution in $(-\infty, 0)$ to the equation

$$\sum_{k=1}^K g_1^{(k)}(z_1) = c \quad (62)$$

i.e., the dominant eigenvalue is the unique parameter in $A^{(k)}(z)$ ($1 \leq k \leq K$) such that the sum of their maximal real eigenvalues is c . ■

The proof closely parallels the proof of Proposition III.4. The main items to note: (58) is satisfied by all eigenvalues; the dominance of $\sum g_1^{(k)}(z)$ for all $z < 0$ as reflected in (61), i.e.,

$$\sum_{k=1}^K g_1^{(k)}(z) \geq \sum_{k=1}^K \text{Re } g_{i(k)}^{(k)}(z) \quad (z < 0) \quad (63)$$

with equality holding only if $i(k) = 1$ for all k ; the aforementioned monotonicity of $\sum g_1^{(k)}(z)$ and range $[\overline{\lambda}, \widehat{\lambda}]$ for $z \in [-\infty, 0]$.

The analog of Proposition III.5 also holds and has a similar proof: the dominant eigenvalue z_1 is monotonic, strictly decreasing with increasing c for $c \in (\overline{\lambda}, \widehat{\lambda})$.

D. Asymptotics

The asymptotic regime is specified by the scaling (37) in which the buffer size $B \rightarrow \infty$ and the buffer overflow probability $p \rightarrow 0$ in a manner parameterized by $\zeta \in [-\infty, 0]$. The following characterizes K sources which satisfy the admission criterion in this asymptotic regime.

Proposition IV.3 Suppose there are K sources $(\mathbf{M}^{(k)}, \boldsymbol{\lambda}^{(k)})$ ($1 \leq k \leq K$). Let the admission criterion be $G(B) \leq p$. Suppose $B \rightarrow \infty$ and $p \rightarrow 0$ in such a manner that $(\log p/B) \rightarrow \zeta \in [-\infty, 0]$.

If $\sum_k g_1^{(k)}(\zeta) < c$, then the admission criterion is satisfied.

If $\sum_k g_1^{(k)}(\zeta) > c$, then the admission criterion is violated.

Here $g_1^{(k)}(\zeta)$ is the maximal real eigenvalue of $A^{(k)}(\zeta) = \left[A^{(k)} - \frac{1}{\zeta} \mathbf{M}^{(k)} \right]$.

Proof: From (9),

$$G(B) = \sum_{i \geq 1} a_i (\boldsymbol{\phi}_i, \mathbf{1}) e^{z_i B}.$$

Here, z_1 is the dominant eigenvalue, so that $z_1 > \text{Re } z_i$ for all $i > 1$ and [see (62)] $\sum g_1^{(k)}(z_1) = c$. Hence,

$$\frac{G(B)}{p} = a_1 (\boldsymbol{\phi}_1, \mathbf{1}) e^{(z_1 - \zeta)B} [1 + o(1)] \quad \text{as } B \rightarrow \infty. \quad (64)$$

Now, from Proposition III.3, $\sum g_1^{(k)}(z)$ decreases as z increases; hence, if $\sum g_1^{(k)}(\zeta) < c$ then $z_1 < \zeta$ and $\{G(B)/p\} \rightarrow 0$ as $B \rightarrow \infty$ and the admission criterion is satisfied. Similarly, if $\sum g_1^{(k)}(\zeta) > c$ then $z_1 > \zeta$ and $\{G(B)/p\} \rightarrow \infty$ and the admission criterion is violated. ■

Now consider the implications of Proposition IV.3 on the admission control problem in which there are, say, J classes of sources. Every source of class j ($1 \leq j \leq J$) is characterized by $(\mathbf{M}^{(j)}, \boldsymbol{\lambda}^{(j)})$. The problem is one of determining the set of all $\mathbf{K} = (K^{(1)}, K^{(2)}, \dots, K^{(J)})$ for which the admission criterion $G_{\mathbf{K}}(B) \leq p$ is satisfied, where $K^{(j)}$ is the number of sources of class j admitted to the multiplexing system.

Corollary.. Let $\mathcal{A}(B, p) = \{\mathbf{K} : G_{\mathbf{K}}(B) \leq p\}$. Also let

$$\mathcal{A} \triangleq \left\{ \mathbf{K} : \sum_j g_1^{(j)}(\zeta) K^{(j)} < c \right\}$$

$$\overline{\mathcal{A}} \triangleq \left\{ \mathbf{K} : \sum_j g_1^{(j)}(\zeta) K^{(j)} \leq c \right\}$$

where $g_1^{(j)}(\zeta)$ is the maximal real eigenvalue of $\left[A^{(j)} - \frac{1}{\zeta} \mathbf{M}^{(j)} \right]$. Then $\mathcal{A} \subseteq \mathcal{A}(B, p) \subseteq \overline{\mathcal{A}}$. ■

In applications of these asymptotic results, we approximate $\mathcal{A}(B, p) \approx \{\mathbf{K} : \sum g_1^{(j)}(\zeta) K^{(j)} < c\}$. Then, except for effects due to the integrality of \mathbf{K} , the acceptance set in \mathbf{K} space is a simplex. The goodness of this approximation is the subject of numerical investigations in the next section.

V. NUMERICAL STUDIES

In the previous sections, we showed that the effective bandwidth of a source is a clearly defined and easily computed quantity. In this section, we numerically investigate three main issues:

- 1) The accuracy of the effective bandwidth when used in admission control and its sensitivity to source burstiness.
- 2) The effect of the variability of the on and off periods on the effective bandwidth.
- 3) The function of the Leaky Bucket regulator as a bandwidth-reducing device.

Figs. 2–4 address issue 1. Fig. 2 displays the boundaries of the acceptance region for two source classes as computed from exact analysis and from using the effective bandwidth approximation. Similar plots were obtained by Gibbens and Hunt [20]. The sources of both classes are on/off with exponentially distributed on and off periods. See (32) for an explanation of the source parameters, which are as follows:

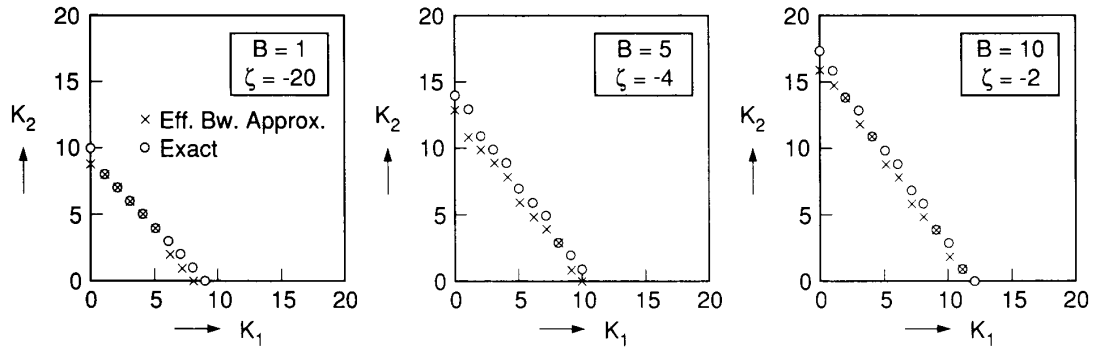


Fig. 2. The acceptance set for two classes of on/off sources with $p = 2.06 \times 10^{-9}$.

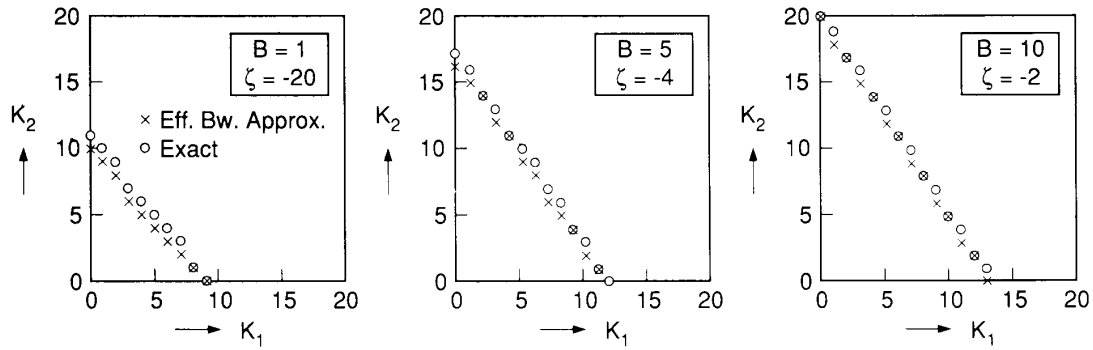


Fig. 3. The acceptance set when the mean on/off periods of both classes are half of those in Fig. 1.

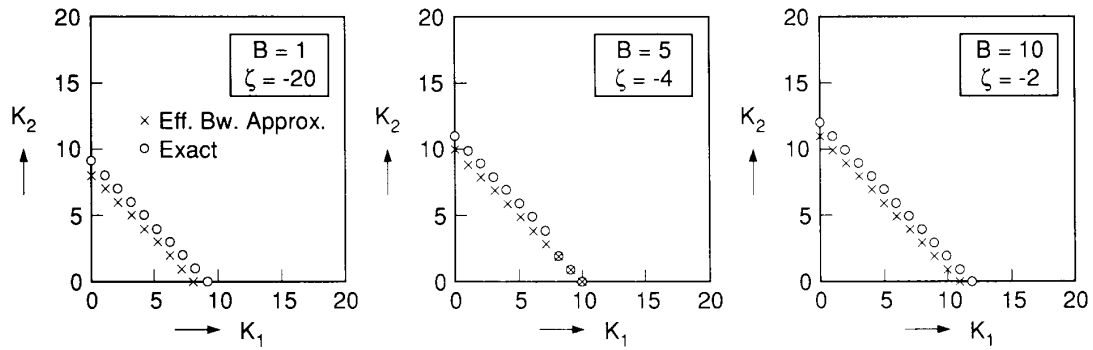


Fig. 4. Acceptance set for two classes of sources with equal effective bandwidths and different mean and peak rates: $p = 2.06 \times 10^{-9}$.

channel capacity $c = 8.43$;

	α	β	λ_1	$\lambda_2 = \hat{\lambda}$
class 1	1.0	1.0	0	1.0
class 2	1.0	2.0	0	1.0

Hence, the on periods of the second class are only half as long. The buffer overflow probability $p = 2.06 \times 10^{-9}$. The buffer size B is varied: $B = 1, 5, 10$ which gives $\zeta = -20, -4$, and -2 , respectively. We point out that these plots differ from those in [20] in that p is constant, while in [20] p is varied with B to keep ζ constant. Also, in Figures 2–4 the data points

are obtained by calculating the maximum acceptable value of K_2 for each value of K_1 .

The reader should observe in these figures that the effective bandwidth approximation provides a *conservative bound* on the acceptance set.

In Fig. 3, the jitteriness of the sources is doubled, i.e., their mean on and off periods are halved. Channel capacity, $c = 8.43$, as in Fig. 2. The effective bandwidth of each source decreases, resulting in an increase in the acceptance set. In Fig. 4, the two source classes have different mean and peak rates but the same effective bandwidth. The parameters are the same as for Fig. 2, except that for class 2 the peak rate $\hat{\lambda} = 1.05$,

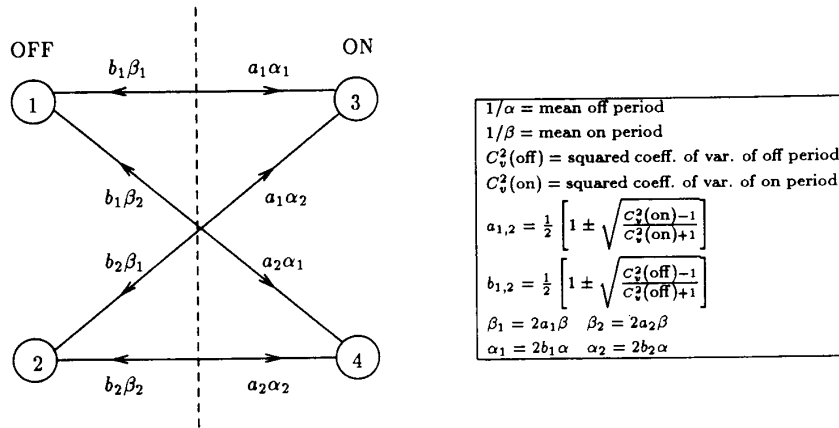


Fig. 5. Four-state source model.

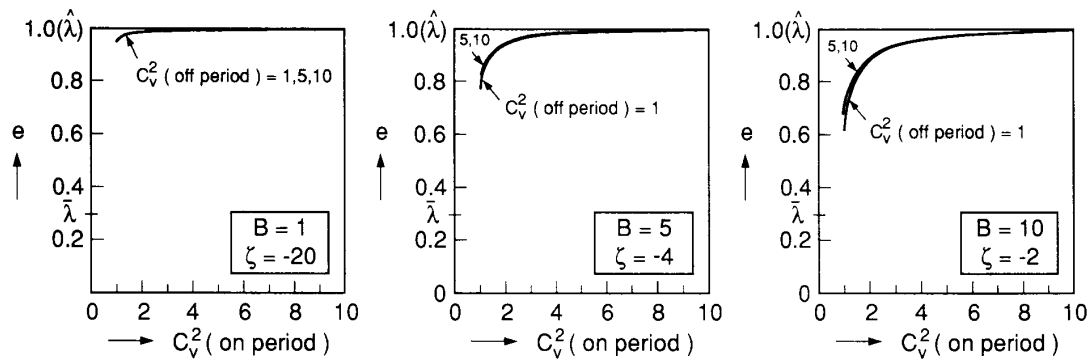


Fig. 6. Effect of C_v^2 on on and off periods. $\alpha = 0.4, \beta = 1, \hat{\lambda} = 1; p = 2.06 \times 10^{-9}$.

1.19, 1.29 for $B = 1, 5, 10$, respectively. The symmetry in the plots confirms that effective bandwidth provides an effective basis for admission control, while mean and peak source rates by themselves do not.

In the previous figures, the on and off periods were assumed to be exponentially distributed. It is of interest to investigate the dependence of the effective bandwidth on the variability of the on and off periods. To generate distributions with a squared coefficient of variations larger than 1, we have chosen the hyperexponential distribution with balanced means. We consider a four-state source model [see Fig. 5], where states 1 and 2 correspond to the off period and states 3 and 4 to the on period. This model allows us (see the equations accompanying the figure) to vary the squared coefficient of variation of the off and on periods, $C_v^2(\text{off})$ and $C_v^2(\text{on})$, while keeping their means constant. In Fig. 6, the effective bandwidth of a source having parameters $\alpha = 0.4, \beta = 1$, and $\hat{\lambda} = 1$ is plotted as a function of $C_v^2(\text{on})$ for various values of $C_v^2(\text{off})$. We observe that the effective bandwidth is sensitive to $C_v^2(\text{on})$ and is less sensitive to $C_v^2(\text{off})$. Fig. 7 displays similar behavior for a source with shorter on and off periods.

In the context of rate-based congestion control, call admission is necessarily complemented by traffic monitoring and regulation. The Leaky Bucket device can act as a traffic policer

as well as a traffic shaper [37],[4],[11],[12],[1]. We consider the simplest form of the Leaky Bucket, which consists of a token pool of size B_T supplied with tokens at rate r . In the model at hand, if an arriving cell finds the token buffer empty, it is marked, allowed into the network, and treated thereafter as a low priority cell. We now examine the effect of the Leaky Bucket on the effective bandwidth of a two-state on/off source. To apply the results derived in this paper, we model the output stream of unmarked, i.e., high priority, cells leaving the Leaky Bucket as a three-state Markov-modulated source as depicted in Fig. 8 (see [11] for a detailed derivation). Fig. 9 plots effective bandwidth versus B_T for different values of r and illustrates the bandwidth-reducing property of the Leaky Bucket. We let the unit of time be the mean length of the on period and the unit of information be the amount generated by the source during an average on period. Thus, the source peak rate and mean rate are equal to 1 and 0.286 units of information per unit of time, respectively. It is seen that the effective bandwidth decreases from a maximum value equal to the source's original effective bandwidth to a minimum value as B_T is decreased from five units of information to 0. The reduction in effective bandwidth is, alternatively, due to the increase in marking probability P_M which increases as B_T decreases, as shown in Fig. 9.

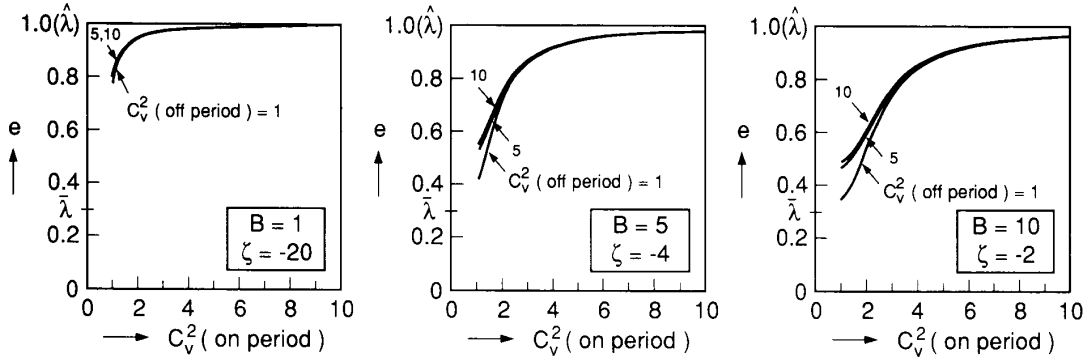


Fig. 7. Effect of C_v^2 on on and off periods. $\alpha = 2, \beta = 5, \hat{\lambda} = 1; p = 2.06 \times 10^{-9}$.

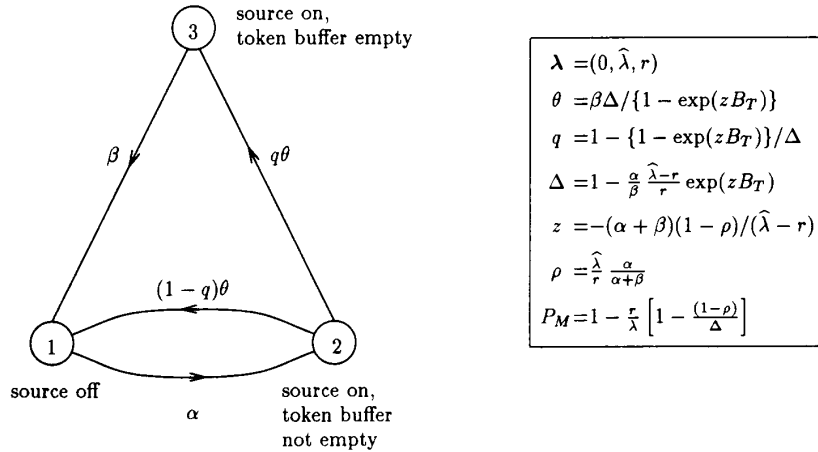


Fig. 8. Approximate Markovian characterization of unmarked cell stream from Leaky Bucket regulator.

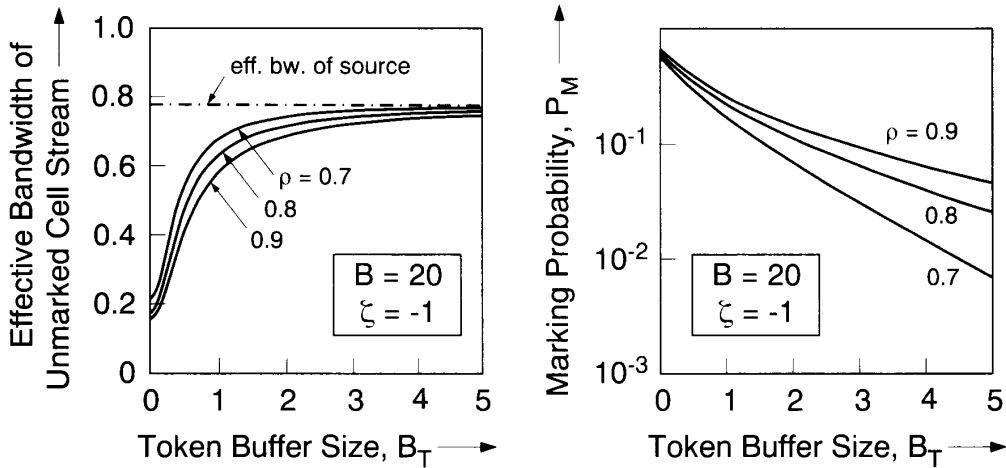


Fig. 9. Effect of the Leaky Bucket regulator on effective bandwidth. $\alpha = 0.1, \beta = 0.25, \hat{\lambda} = 1; p = 2.06 \times 10^{-9}$. Token rate $r = 0.41, 0.35, 0.32$ for regulator's $\rho = 0.7, 0.8, 0.9$, respectively.

VI. EFFECTIVE BANDWIDTH OF A VIDEO TELECONFERENCING TRAFFIC SOURCE

In this section, we demonstrate that a realistic traffic source derived from measurements has an effective bandwidth which

is easy to calculate. The model, which is due to Heyman *et al.* [21], is for traffic from video teleconferencing services such as would be provided over ATM-based networks. Beginning with a 30-minute sequence of video teleconference data, the

authors fit a variety of autoregressive and Markov chain models and conclude that the only model sufficiently accurate for use in traffic studies is a multistate Markov chain model. A version of this model which they recommend is the discrete autoregressive model DAR(1). While this model is structurally simple, the number of states is about 600 in the version which is used in their simulation experiments. We consider the continuous-time version of the DAR(1) model and show that its effective bandwidth is readily calculated, even when the number of states is large.

In our model, the infinitesimal generator

$$\mathbf{M} = \alpha[-\mathbf{I} + \mathbf{1}\langle\mathbf{f}] \quad (65)$$

where $\langle\cdot\rangle$ denotes the outer product of vectors, and in this case $\mathbf{1}\langle\mathbf{f}$ is a matrix in which every row is identical to \mathbf{f} . The form in (65) is obtained by identifying the transition matrix $\mathbf{P} = [\rho\mathbf{I} + (1-\rho)\mathbf{1}\langle\mathbf{f}]$ in [21] with $\exp(\mathbf{M}\Delta)$, where ρ is a first-order autocorrelation coefficient, \mathbf{f} is a probability vector, i.e., $\mathbf{f} \geq \mathbf{0}$, $\langle\mathbf{f}, \mathbf{1}\rangle = 1$, which is derived from the negative binomial distribution, and Δ is a small parameter associated with the time discretization. Since $\mathbf{P} = \exp(\mathbf{M}\Delta) \approx \mathbf{I} + \mathbf{M}\Delta$, (65) is obtained on identifying α with $(1-\rho)/\Delta$. The model in [21] implies a rate vector λ for our source with a linear dependence of λ_i on i . However, it will be equally convenient to let λ be arbitrary.

The effective bandwidth of the source, $e(\mathbf{M}, \lambda; B, p)$, is the maximal real solution e of the equation

$$\left| \mathbf{A} - \frac{1}{\zeta} \mathbf{M} - e\mathbf{I} \right| = 0 \quad (66)$$

where $\zeta = \log p/B$. To evaluate the determinant, we make use of the identity

$$|\mathbf{A} - \mathbf{a}\langle\mathbf{b}| = |\mathbf{A}|(1 - \langle\mathbf{a}\mathbf{A}^{-1}, \mathbf{b}\rangle) \quad (67)$$

and obtain

$$\left| \mathbf{A} - \frac{1}{\zeta} \mathbf{M} - e\mathbf{I} \right| = \left[\prod_i \left(\lambda_i - e + \frac{\alpha}{\zeta} \right) \right] \left[1 - \frac{\alpha}{\zeta} \sum_i \frac{f_i}{\lambda_i - e + \alpha/\zeta} \right]. \quad (68)$$

Hence, the effective bandwidth is the maximal real solution to the equation

$$\frac{\alpha}{\zeta} \sum_i \frac{f_i}{\lambda_i + \alpha/\zeta - e} = 1. \quad (69)$$

Suppose that $\lambda_1 < \lambda_2 < \dots < \lambda_{N-1} < \lambda_N$. The function on the left is monotonic, increasing in each of the intervals $(-\infty, \lambda_1 + \alpha/\zeta)$, $(\lambda_1 + \alpha/\zeta, \lambda_2 + \alpha/\zeta)$, $(\lambda_2 + \alpha/\zeta, \dots, (\lambda_{N-1} + \alpha/\zeta, \lambda_N + \alpha/\zeta)$, $(\lambda_N + \alpha/\zeta, \infty)$; the function approaches $\pm\infty$ as the singularities at $(\lambda_i + \alpha/\zeta)$ are approached from the left and right, and approaches 0 as $e \rightarrow -\infty$ and as $e \rightarrow \infty$. Hence,

Proposition VI.1 The effective bandwidth of source (\mathbf{M}, λ) , where \mathbf{M} is given in (65), is the unique real solution e in the interval $(\lambda_{N-1} + \alpha/\zeta, \lambda_N + \alpha/\zeta)$ to (69). ■

The use of this result in evaluating the effect of smoothing on the statistical multiplexing gain is the subject of current investigations.

VII. EFFECTIVE BANDWIDTH OF MARKOV MODULATED POISSON SOURCES

In this section, we show that there are closely related concepts and arguments that apply to queues of jobs or packets. This parallelism between fluid flow and point processes has been noted before in [15] and [13], in particular for the decompositions of the eigenvalue problem. We shall show here that the parallelism also extends to the inverse eigenvalue problem, the qualitative properties of the maximal real eigenvalue as a function of the parameter of the problem, and the concept of effective bandwidth.

We begin, as in Section III, with a single source (\mathbf{M}, λ) , where \mathbf{M} is the irreducible infinitesimal generator of a controlling Markov chain. The source emits packets in a Poisson stream at rate λ_s when in state s ($s \in \mathcal{S}$). Let $\mathbf{A} = \text{diag}(\lambda)$. The packet length is exponentially distributed. The server is the output channel to the multiplexing buffer and has constant capacity or rate. The rate parameter μ is the ratio of the channel capacity to the mean packet length. The vector \mathbf{w} , the mean source rate $\bar{\lambda}$, and the peak rate $\bar{\lambda}$ are all defined as in Section II-A. The ergodicity condition is $\bar{\lambda} < \mu$.

Let the stationary state distribution of the multiplexing system be denoted by $\mathbf{p}(n) = \{p_s(n) | s \in \mathcal{S}\}$, where

$$p_s(n) = \Pr \left(\sum = s, X = n \right) \quad (s \in \mathcal{S}; n = 0, 1, \dots). \quad (70)$$

The balance equations are

$$\begin{aligned} \mathbf{0} &= \mathbf{p}(n)[\mathbf{M} - \mathbf{A}] + \mu\mathbf{p}(n+1) & (n=0) \\ &= \mathbf{p}(n-1)\mathbf{A} + \mathbf{p}(n)[\mathbf{M} - \mathbf{A} - \mu\mathbf{I}] + \mu\mathbf{p}(n+1) & (n \geq 1). \end{aligned} \quad (71)$$

The spectral representation of the solution to the balance equations

$$\mathbf{p}(n) = \sum_{i: |z_i| < 1} a_i \phi_i z_i^n \quad (n \geq 0) \quad (72)$$

where (z_i, ϕ_i) is an eigenvalue/eigenvector pair satisfying the eigenvalue equation

$$\phi[\mu z^2 \mathbf{I} + z(\mathbf{M} - \mathbf{A} - \mu\mathbf{I}) + \mathbf{A}] = \mathbf{0}. \quad (73)$$

Let the eigenvalues with modulus less than unity be indexed so

$$1 > |z_1| \geq |z_2| \geq \dots \quad (74)$$

The coefficients $\{a_i\}$ in (72) are obtained from the normalization conditions $\sum \mathbf{p}(n) = \mathbf{w}$. On substituting the solution in (72), we obtain

$$\mathbf{p}(n) = \mathbf{w}[\mathbf{I} - \mathbf{R}] \mathbf{R}^n \quad (n \geq 0) \quad (75)$$

where

$$\mathbf{R} \triangleq \Phi^{-1} \mathbf{Z} \Phi, \quad (76)$$

$\mathbf{Z} = \text{diag}\{z_1, z_2, \dots\}$ and Φ is the matrix with rows ϕ_1, ϕ_2, \dots . From (75),

$$G(n) \triangleq \Pr(X \geq n) = \langle \mathbf{w} \mathbf{R}^n, \mathbf{1} \rangle \quad (n \geq 0). \quad (77)$$

Equation (75) is also the well-known matrix-geometric form due to Neuts [31], who has shown that the *rate matrix* \mathbf{R} has spectral radius less than unity and is the minimal nonnegative solution to a matrix quadratic equation. Hence, \mathbf{R} has a Perron root, the eigenvalue of maximum modulus which is real, simple, and in $(0, 1)$. From the spectral expansion of \mathbf{R} in (76), we infer that the Perron root is z_1 , the system eigenvalue [see (73) and (74)]. Hence, z_1 is real, simple, and in $(0, 1)$. Since $z_1 > |z_i|$ for all $i > 1$, z_1 is called the dominant eigenvalue.

The inverse eigenvalue problem is obtained as in Section II-B. On writing $g(z) = \mu$,

$$g(z)\phi = \phi \mathbf{A}(z) \quad (78)$$

where

$$\mathbf{A}(z) \triangleq \frac{1}{z} \mathbf{A} + \frac{1}{1-z} \mathbf{M}. \quad (79)$$

$\mathbf{A}(z)$ is irreducible and essentially nonnegative for $z \in (0, 1)$. Let $g_1(z)$ be its maximal real eigenvalue. We have proven

Proposition VII.1

- 1) $g_1(z) \sim \hat{\lambda}/z$ as z approaches 0 from the right, and $g_1(1) = \bar{\lambda}$.
- 2) $g_1'(z) < 0$ for $z \in (0, 1)$.
- 3) The dominant eigenvalue z_1 is the unique solution in $(0, 1)$ satisfying $g_1(z_1) = \mu$. ■

We next examine the admission criterion $\{G(B) \leq p\}$ in the asymptotic regime of large buffers B and small overflow probabilities p . Our result is

Proposition VII.2 Suppose $B \rightarrow \infty$ and $p \rightarrow 0$ in such a manner that $(\log p/B) \rightarrow \zeta \in [-\infty, 0]$.

If $g_1(e^\zeta) < \mu$, then the admission criterion is satisfied;

If $g_1(e^\zeta) > \mu$, then the admission criterion is violated where $g_1(e^\zeta)$ is the maximal real eigenvalue of $\mathbf{A}(e^\zeta) = \frac{1}{e^\zeta} \mathbf{A} + \frac{1}{1-e^\zeta} \mathbf{M}$. ■

On the basis of this result, the effective bandwidth of the single source

$$e(\mathbf{M}, \boldsymbol{\lambda}; B, p) = g_1(e^\zeta). \quad (80)$$

For a two-state MMPP source with $(\mathbf{M}, \boldsymbol{\lambda})$ defined in (32),

$$g_1(z) = \frac{1}{2} \left[\frac{\lambda_1 + \lambda_2}{z} - \frac{\alpha + \beta}{1-z} \right] - \frac{1}{2} \sqrt{\left(\frac{\lambda_1 + \lambda_2}{z} - \frac{\alpha + \beta}{1-z} \right)^2 - 4 \left(\frac{\lambda_1 \lambda_2}{z^2} - \frac{\alpha \lambda_2 + \beta \lambda_1}{z(1-z)} \right)}. \quad (81)$$

We next investigate, as in Section IV, the decomposition of the expression in (80) when the source $(\mathbf{M}, \boldsymbol{\lambda})$ is the aggregate of K sources $(\mathbf{M}^{(k)}, \boldsymbol{\lambda}^{(k)})$ ($1 \leq k \leq K$). The key coupled eigenvalue problem in (58) carries over, with $\mathbf{A}^{(k)}(z) \triangleq \frac{1}{z} \mathbf{A}^{(k)} + \frac{1}{1-z} \mathbf{M}^{(k)}$. With the benefit of Proposition VII.1, we have proven

Proposition VII.3 Suppose there are K sources $(\mathbf{M}^{(k)}, \boldsymbol{\lambda}^{(k)})$ ($1 \leq k \leq K$) supplying the multiplexing buffer. Let the admission criterion and asymptotic regime be as in Proposition VII.2.

If $\sum_k g_1^{(k)}(e^\zeta) < \mu$, then the admission criterion is satisfied;

If $\sum_k g_1^{(k)}(e^\zeta) > \mu$, then the admission criterion is violated where $g_1^{(k)}(e^\zeta)$ is the maximal real eigenvalue of $\mathbf{A}^{(k)}(e^\zeta)$. ■

TABLE I
THE NUMBER OF ADMISSIBLE SOURCES OBTAINED BY THE EFFECTIVE BANDWIDTH APPROXIMATION AND EXACT CALCULATIONS.

μ	Case 1		Case 2	
	K_e	K^*	K_e	K^*
50	12	12	10	11
100	24	24	21	22
150	36	36	32	33
200	48	49	43	44
250	60	61	54	56
300	73	74	65	67

Thus, the simple additive structure to the effective bandwidth of K sources exists in both the fluid and queueing frameworks. All the simplifying consequences discussed in Section IV-D carry over.

We give numerical results on the use of the effective bandwidth approximation to the admission control of a single class of two-state MMPP sources. We consider two cases: Case 1 and Case 2. The sources in the two cases have the same mean rate ($\bar{\lambda} = 3.33$) and different burstiness characteristics. In both cases, $B = 200$ and $p = 10^{-7}$, which gives $e^\zeta = 0.9226$. The source parameters and effective bandwidth $e = e(\mathbf{M}, \boldsymbol{\lambda}; B, p)$ of the sources as computed by (80) and (81) are given below.

	α	β	λ_1	$\lambda_2 = \hat{\lambda}$	e
Case 1	1	1.	0	6.667	4.110
Case 2	1	4.	0	16.667	4.568

In Table I, we compare the admissible number of sources K_e , computed by the effective bandwidth approximation, with K^* , obtained by exact calculations [15]. The comparison is carried out for a range of values of μ , the service rate.

The main observation in Table I is that K^* is tightly and conservatively bounded by K_e for sufficiently large B . Also, the admissible number in Case 1 is consistently larger than that of Case 2. This is because sources of Case 2 are more bursty.

In recent work [14], we have extended the results of this section to phase renewal (PH renewal) processes [31].

APPENDIX KRONECKER ALGEBRA

The Kronecker product $\mathbf{A} \otimes \mathbf{B}$ of the matrix \mathbf{A} of dimension $p \times q$ and the matrix \mathbf{B} of dimension $m \times n$ is the matrix of dimension $pm \times qn$ obtained by replacing each element a_{ij} of the matrix \mathbf{A} by the full matrix $a_{ij} \mathbf{B}$. (See, for example, [3].)

The Kronecker sum of $\mathbf{A}(n \times n)$ and $\mathbf{B}(m \times m)$ denoted by $\mathbf{A} \oplus \mathbf{B}$ is defined as

$$\mathbf{A} \oplus \mathbf{B} = \mathbf{A} \otimes \mathbf{I}_m + \mathbf{I}_n \otimes \mathbf{B}$$

where \mathbf{I}_m and \mathbf{I}_n are the identity matrices of order m and n , respectively. The operation \otimes is associative but not commutative, and the same holds true for \oplus .

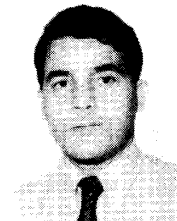
The following properties, which are proven in [6], [3], and [18], are used in this paper. All matrices (vectors) are assumed to be of appropriate order.

- 1) $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{AC}) \otimes (\mathbf{BD})$.

- 2) Let $\lambda_1, \lambda_2, \dots, \lambda_n$ be the eigenvalues of the matrix A with corresponding eigenvectors $\alpha_1, \alpha_2, \dots, \alpha_n$ and let $\mu_1, \mu_2, \dots, \mu_m$ be the eigenvalues of B with corresponding eigenvectors $\beta_1, \beta_2, \dots, \beta_m$. Then, the eigenvalues of $A \oplus B$ are the nm sums $\lambda_i + \mu_j$ with corresponding eigenvectors $\alpha_i \otimes \beta_j$, $i = 1, 2, \dots, n$ and $j = 1, 2, \dots, m$.

REFERENCES

- [1] V. Anantharam and P. Konstantopoulos, "Burst reduction properties of the leaky bucket in ATM networks," preprint, 1992.
- [2] D. Anick, D. Mitra, and M. M. Sondhi, "Stochastic theory of a data handling system with multiple sources," *Bell Syst. Tech. J.*, vol. 61, pp. 1871-1894, 1982.
- [3] R. Bellman, *Introduction to Matrix Analysis*, 2nd ed. New York, NY: McGraw-Hill, 1970.
- [4] A. Berger, "Performance analysis of a rate-control throttle where tokens and jobs queue," *IEEE J. Select. Areas Commun.*, vol. 9, pp. 165-170, 1991.
- [5] B. Bensaou, J. Guibert, J. W. Roberts, and A. Simonian, "Performance of an ATM multiplexer queue in the fluid approximation using the Benes approach," preprint, 1992.
- [6] J. W. Brewer, "Kronecker products and matrix calculus in systems theory," *IEEE Trans. Circ. Syst.*, vol. 25, pp. 772-781, 1978.
- [7] C.-S. Chang, "Stability, queue length and delay, part II: Stochastic queueing networks," preprint, 1992.
- [8] J. E. Cohen, "Random evolutions and the spectral radius of a non-negative matrix," *Math. Proc. Camb. Phil. Soc.*, vol. 86, pp. 345-350, 1979.
- [9] J. E. Cohen, "Convexity of the dominant eigenvalue of an essentially nonnegative matrix," in *Proc. AMS '81*, pp. 657-658.
- [10] G. Debreau and I. N. Herstein, "Nonnegative square matrices," *Econometrica*, vol. 21, pp. 597-607, 1953.
- [11] A. I. Elwalid and D. Mitra, "Analysis and design of rate-based congestion control of high speed networks, part I: Stochastic fluid models, access regulation," *Queueing Syst.*, vol. 9, pp. 29-64, 1991.
- [12] 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 '92*, Florence, Italy, pp. 415-425.
- [13] A. I. Elwalid and D. Mitra, "Markovian arrival and service communication systems: Spectral expansions, separability, Kronecker-product forms," preprint, 1993.
- [14] A. I. Elwalid and D. Mitra, "Effective bandwidth of bursty, variable rate sources for admission control to B-ISDN," in *Proc. ICC '93*, Geneva, Switzerland, pp. 1325-1330.
- [15] A. I. Elwalid, D. Mitra, and T. E. Stern, "Statistical multiplexing of Markov-modulated sources: Theory and computational algorithms," in *Teletraffic and Data Traffic in a Period of Change*, 1991, pp. 495-500.
- [16] S. Friedland, "Convex spectral functions," *Linear and Multilinear Algebra*, vol. 9, pp. 299-316, 1981.
- [17] F. R. Gantmacher, *The Theory of Matrices*. New York, NY: Chelsea, 1960, vol. 2.
- [18] A. Graham, *Kronecker Products and Matrix Calculus with Applications*. Chichester: Ellis Horwood, 1981.
- [19] 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, pp. 968-981, 1991.
- [20] R. J. Gibbens and P. J. Hunt, "Effective bandwidths for the multi-type UAS channel," *Queueing Syst.*, vol. 9, pp. 17-28, 1991.
- [21] D. Heyman, A. Tabatabai, and T. V. Lakshman, "Statistical analysis and simulation study of video teleconference traffic in ATM network," *IEEE Trans. Circ. Syst. for Video Technol.*, vol. 2, pp. 49-59, 1992.
- [22] J. Y. Hui, "Resource allocation for broadband networks," *IEEE J. Select. Areas Commun.*, vol. 6, pp. 1598-1608, 1988.
- [23] J. F. C. Kingman, "A convexity property of positive matrices," *Quart. J. Math. Oxford Ser.*, vol. 12, pp. 283-284, 1961.
- [24] F. P. Kelly, "Effective bandwidths at multi-type queues," *Queueing Syst.*, vol. 9, pp. 5-15, 1991.
- [25] L. Kosten, "Stochastic theory of data-handling systems with groups of multiple sources," in *Performance of Computer Communication Systems*, H. Rudin and W. Bux, Eds. Amsterdam, The Netherlands: Elsevier, 1984, pp. 321-331.
- [26] L. Kosten, "Liquid models for a type of information buffer problem," *Delft Prog. Rep.* 11, 1986, pp. 71-86.
- [27] G. Kesidis and J. Walrand, "Effective bandwidths for multiclass Markov fluid and other ATM sources," U.C. Berkeley rep., 1992.
- [28] K. Lindberger, "Analytical methods for the traffic problems with statistical multiplexing in ATM-networks," in *Teletraffic and Data Traffic in a Period of Change*, 1991.
- [29] B. Maglaris, P. Anastassiou, P. Sen, G. Karlsson, and J. D. Robbins, "Performance models of statistical multiplexing in packet video communications," *IEEE Trans. Commun.*, vol. 36, pp. 834-843, 1988.
- [30] D. Mitra, "Stochastic theory of a fluid model of producers and consumers coupled by a buffer," *Adv. Appl. Prob.*, vol. 20, pp. 646-676, 1988.
- [31] M. F. Neuts, *Matrix-Geometric Solutions to Stochastic Models*. Baltimore, MA: John Hopkins Univ. Press, 1981.
- [32] W. H. Press, B. P. Flannery, S. A. Teukolsky, and W. T. Vetterling, *Numerical Recipes in C*. Cambridge, MA: Cambridge Univ. Press, 1988.
- [33] J. W. Roberts, "Performance evaluation and design of multiservice networks," Final Report of the COST 224 Project, Commission of the European Communities, 1992.
- [34] J. W. Roberts, "Traffic control in the B-ISDN," to appear in *Comput. Netw. and ISDN Syst.*
- [35] T. E. Stern and A. I. Elwalid, "Analysis of a separable Markov-modulated rate model for information-handling systems," *Adv. Appl. Prob.*, vol. 23, pp. 105-139, 1991.
- [36] E. Seneta, *Nonnegative Matrices: An Introduction to Theory and Applications*. London: Allen and Unwin, 1973.
- [37] M. Sidi, W. Z. Liu, I. Cidon, and I. Gopal, "Congestion control through input rate regulation," in *Proc. GLOBECOM '89*, Dallas, TX, pp. 1764-1768.
- [38] K. Sohraby, "On the asymptotic behavior of heterogeneous statistical multiplexer with applications," in *Proc. IEEE INFOCOM '92*, pp. 839-847.
- [39] J. H. Wilkinson, *The Algebraic Eigenvalue Problem*. Oxford, U.K.: Clarendon, 1965.
- [40] W. Whitt, "Tail probabilities with statistical multiplexing and effective bandwidths for multi-class queues," to appear in *Telecommun. Syst.*, 1992.



Anwar I. Elwalid (M'91) received the B.S. degree in electrical engineering from the Polytechnic Institute of New York, Brooklyn, and the M.S. and Ph.D. degrees in electrical engineering from Columbia University, New York.

Since 1991, he has been with the Mathematics of Networks and Systems Research Department at Bell Laboratories, Murray Hill, NJ. His research interests are in communication and computer networks, queueing, and stochastic systems. He is a member of Tau Beta Pi (National Engineering Honor Society) and Sigma Xi (Scientific Research Society).

Since 1991, he has been with the Mathematics of Networks and Systems Research Department at Bell Laboratories, Murray Hill, NJ. His research interests are in communication and computer networks, queueing, and stochastic systems. He is a member of Tau Beta Pi (National Engineering Honor Society) and Sigma Xi (Scientific Research Society).



Debasis Mitra (M'75-SM'82-F'89) was born in India in 1944. He received the B.Sc. and Ph.D. degrees in electrical engineering from London University in 1964 and 1967, respectively.

He joined Bell Laboratories as a Member of Technical Staff in 1968. Since 1986, he has been Head of the Mathematics of Networks and Systems Research Department. During the fall semester of 1984, he was Visiting McKay Professor at the University of California, Berkeley. He has been involved in the development of asymptotic theories for large queueing networks and their incorporation in the software package PANACEA. He has worked on asynchronous computations for parallel processors, kanban manufacturing models, and state-dependent routing on circuit-switched networks. In the late 1970's and early 1980's, he worked on the analyses of statistical multiplexing of bursty traffic by means of stochastic fluid models. More recently, he has returned to study the control and design of networks subject to bursty traffic. He has also been involved in a fundamental investigation of the role of feedback in high-speed wide area networks. He is the recipient of awards given by the Institution of Electrical Engineers, United Kingdom, the Bell System Technical Journal, and of the Steven O. Rice Prize Paper Award and the Guillemin-Cauer Prize Paper Award of the IEEE.

Dr. Mitra is a member of the ACM, SIAM, ORSA, and IFIP Working Group 7.3.