

**Analysis of Separable Markov-Modulated Rate Models for
Information-Handling Systems**



Thomas E. Stern; Anwar I. Elwalid

Advances in Applied Probability, Vol. 23, No. 1. (Mar., 1991), pp. 105-139.

Stable URL:

<http://links.jstor.org/sici?sici=0001-8678%28199103%2923%3A1%3C105%3AAOSMRM%3E2.0.CO%3B2-3>

Advances in Applied Probability is currently published by Applied Probability Trust.

Your use of the JSTOR archive indicates your acceptance of JSTOR's Terms and Conditions of Use, available at <http://www.jstor.org/about/terms.html>. JSTOR's Terms and Conditions of Use provides, in part, that unless you have obtained prior permission, you may not download an entire issue of a journal or multiple copies of articles, and you may use content in the JSTOR archive only for your personal, non-commercial use.

Please contact the publisher regarding any further use of this work. Publisher contact information may be obtained at <http://www.jstor.org/journals/apt.html>.

Each copy of any part of a JSTOR transmission must contain the same copyright notice that appears on the screen or printed page of such transmission.

JSTOR is an independent not-for-profit organization dedicated to creating and preserving a digital archive of scholarly journals. For more information regarding JSTOR, please contact support@jstor.org.

ANALYSIS OF SEPARABLE MARKOV-MODULATED RATE MODELS FOR INFORMATION-HANDLING SYSTEMS

THOMAS E. STERN* AND

ANWAR I. ELWALID, **Columbia University*

Abstract

In many communication and computer systems, information arrives to a multiplexer, switch or information processor at a rate which fluctuates randomly, often with a high degree of correlation in time. The information is buffered for service (the server typically being a communication channel or processing unit) and the service rate may also vary randomly. Accurate capture of the statistical properties of these fluctuations is facilitated by modeling the arrival and service rates as superpositions of a number of independent finite state reversible Markov processes. We call such models separable Markov-modulated rate processes (MMRP).

In this work a general mathematical model for separable MMRPs is presented, focusing on Markov-modulated continuous flow models. An efficient procedure for analyzing their performance is derived. It is shown that the 'state explosion' problem typical of systems composed of a large number of subsystems, can be circumvented because of the separability property, which permits a decomposition of the equations for the equilibrium probabilities of these systems. The decomposition technique (generalizing a method proposed by Kosten) leads to a solution of the equilibrium equations expressed as a sum of terms in Kronecker product form. A key consequence of decomposition is that the computational complexity of the problem is vastly reduced for large systems. Examples are presented to illustrate the power of the solution technique.

KRONECKER PRODUCT FORM; CONTINUOUS FLOW MODEL

1. Introduction

In many communication and computer systems, information arrives to a multiplexer, switch or information processor at a rate which fluctuates randomly, often with a high degree of correlation in time. The information is buffered for service (the server typically being a communication channel or processing unit) and the service rate may also vary randomly. Accurate capture of the statistical properties of these fluctuations is facilitated by modeling the arrival and service rates as Markov processes. We call such models Markov-modulated rate processes (MMRP). In

Received 7 June 1989; revision received 27 February 1990

* Postal address for both authors: Department of Electrical Engineering and Center for Telecommunications Research, Columbia University, New York, NY 10027, USA.

Part of this work was performed when the first author was on sabbatical leave at INRIA Centre de Sophia Antipolis 06565 Valbonne, France. This work was partially supported by the National Science Foundation under grant CDR-84-21402, and the office of Naval Research under grant N00014-85-k0371.

some applications (typically, realtime communications such as voice or video) the information being processed may be realistically modeled as a Markov-modulated continuous flow process, while in others, (e.g. computer data) a point process model (e.g. Markov-modulated Poisson) may be more appropriate.

In practical design and performance evaluation applications involving MMRPs, it is often necessary to obtain detailed information regarding buffer occupancy distributions. It is frequently the tails of these distributions which contain critical design and performance information. In these cases mean values are useless. The analytical difficulties are compounded by the fact that in realistic models of 'moderate size' the modulating process will often have thousands of states. This 'state explosion' problem poses various computational difficulties: enormous memory requirements, very long computation time, and complete breakdown of numerical algorithms.

The objective of this work is to present an analysis of these systems in as general a setting as possible, using a decomposition technique which yields insight into system behavior, gives complete information on buffer occupancy distributions, and which is computationally feasible for very large systems. We focus on the continuous flow model, which has been termed a fluid flow model by Anick et al. [1] or a process with chain-dependent growth rate by Keilson and Rao [9], [10]. However, the general approach can be extended to point processes as well. The literature abounds with examples of Markov-modulated rate processes. Existing analytical techniques for such systems generally fall into three categories: (1) Exact analysis of systems of relatively small size, Fischer [6] or with a special structure [22], [21], [18], [1], [15], [13], [14], [16], (2) analysis using simplifying approximations (typically approximating the modulating process as uncorrelated), and (3) asymptotic approaches, e.g. [19], [4], [24]. Of the work just cited, [1] and [13] are of special interest here. Anick et al. [1], showed that when the MMRP consists of the superposition of a finite number of independent identical indistinguishable continuous flow sources each modulated by a two-state Markov process, and the service capacity is constant, it is possible to derive closed-form solutions for the buffer-occupancy distributions, as well as very simple asymptotic approximations. They exhibited numerical results for systems with up to 767 states. Kosten [13] showed how [1] could be extended to systems composed of a superposition of several independent subsystems, each one of which is of the type considered in [1]. In more recent work [15], carried out independently and simultaneously with ours, Mitra extended [1] to systems with variable service capacity, and obtained certain results on the eigenvalues of reversible systems which are similar to some of ours. Our work was motivated by the idea that the basic approach of [1], [13] should be generalizable to much larger classes of systems. We show herein how the main features of both [1] and [13] carry over to MMRPs in which the underlying Markov process for arrivals and/or service are superpositions of any number of independent non-identical finite state reversible

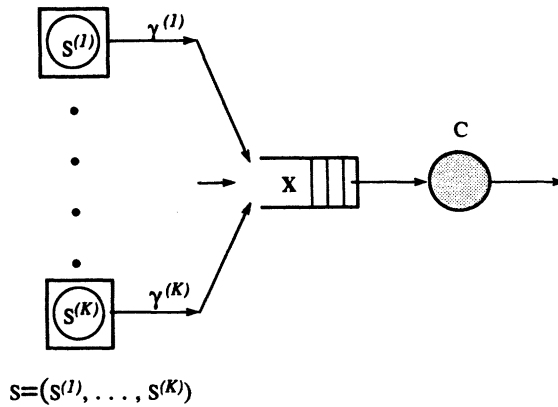


Figure 1. Markov-modulated sources model

processes of arbitrary structure. We call these *separable* MMRPs. To exhibit the relation between this work and its antecedents, we keep our notation consistent with that of [1], [13] wherever possible.

To fix ideas consider two typical problems.

1.1. *Problem 1.* Suppose that K independent information sources are served by a common communication channel of capacity C data units/s, where the data units (du) might be packets, bytes or bits (see Figure 1). The k th source generates information at a rate $\gamma^{(k)}(s^{(k)})$ data units/s, where $s^{(k)}$ is the state of a Markov process which ‘modulates’ the source arrival process. We assume that the data units arrive as a continuous flow while the source remains in a given state. (In another version of the problem the source arrival rate could be a Markov-modulated Poisson process, with γ being the state-dependent arrival rate.) Since the combined arrival rate may sometimes exceed C , temporarily overloading the channel, a buffer is provided, whose content X is a continuous random variable.

The equation for the evolution of the buffer occupancy X is

$$(1) \quad \frac{dX}{dt} = \gamma(S(t)) - C, \quad X > 0$$

where γ is the combined arrival rate from all sources, and $S(t)$ is the combined state $(s^{(1)}, s^{(2)}, \dots, s^{(K)})$, of the K sources at time t . The buffer introduces a random waiting time X/C , which is a critical performance parameter in real time applications. Another important parameter is probability of buffer overflow. Thus, the analysis problem here consists of determining the complete probability distribution of the random variable X , from which these and other performance parameters of interest can be directly determined.

At this point we have not specified the functional form of the arrival rates $\gamma^{(k)}(s^{(k)})$ nor the structure of the underlying Markov process. In a typical

application, all sources might be identical, with each source modulated by an independent N -state process, representing the statistical features of a data terminal [1], or an encoded voice or video source. If $N = 2$ we have the case treated in [1]. Models with $N > 2$ are useful for modeling voice sources whose periods of silence have complex statistics, (for example, short silences between words/syllables and long silences while listening to the other speaker in a dialog). They are also useful as models of encoded video sources which produce information at several different rates, depending on visual activity (movement).

1.2. *Problem 2.* A communication channel of capacity C serves two classes of traffic: x and y (see Figure 2). Class x is managed by delay, that is, a buffer is provided to store all class x information that cannot be immediately served. Class y is managed by loss: any information that cannot be immediately served is discarded. This type of arrangement might be suitable, for example, in a situation where class x is computer data, for which loss cannot be tolerated, but some delay is tolerable, while class y is voice or video traffic, which must be transmitted with minimum delay, but for which some small probability of loss is acceptable. The service strategy is one of partial sharing of the communication channel, with a fraction, $(1 - \alpha)$, reserved for class x and α for class y . If one class does not completely occupy its reserved portion of the channel, the residual capacity is available for the other class. Each traffic class is assumed to arrive as a continuous flow generated by a Markov-modulated source of the type described in Problem 1, with arrival rates $\gamma^{(x)}(s^{(x)})$ and $\gamma^{(y)}(s^{(y)})$ for the x and y sources respectively. The combined system state is $S = (s^{(x)}, s^{(y)})$.

Let X be the class x buffer content and L the accumulated loss of class y information. Then, according to the partial sharing rules indicated above, the equations of evolution of X and L are

$$(2) \quad \frac{dX}{dt} = \gamma^{(x)}(s^{(x)}) + \min(\alpha C, \gamma^{(y)}(s^{(y)})) - C \quad X > 0,$$

$$(3) \quad \frac{dL}{dt} = \max(\gamma^{(y)}(s^{(y)}) - \alpha C, 0) \quad X > 0,$$

$$(4) \quad = \max(\gamma^{(y)}(s^{(y)}) + \gamma^{(x)}(s^{(x)}) - C, 0) \quad X = 0.$$

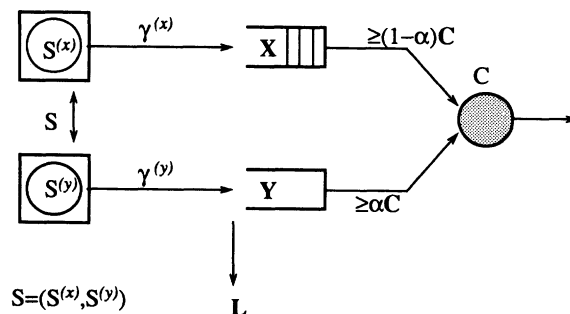


Figure 2. Example of a shared channel

Note from Equation (2) that the channel sharing strategy provides a minimum service rate $(1 - \alpha)C$ to class x traffic at all times, and this rate may increase to the total channel capacity C if the class y arrival rate drops to zero. Whenever the X buffer is non-empty, class y traffic experiences a loss rate equal to the difference between its arrival rate and its reserved service rate, αC . If $X = 0$, (which is only possible if the right-hand side of (2) is ≤ 0), the class y traffic has more capacity available to it, and the loss rate diminishes to zero as the total arrival rate decreases below C . The case $\alpha = 1$ gives complete priority to class y traffic.

The arrival statistics of the two traffic classes in this problem can be quite general. Each of the classes might be modeled as a Markov-modulated source, with $N^{(x)}$ and $N^{(y)}$ states in the underlying processes for sources x and y respectively. If the system is large, it is likely that each source would be composed of the superposition of many, possibly identical smaller sources, each with a relatively small number of states.

As will be shown below, the models in the above problems are both examples of separable MMRPs. In each, the state explosion problem is evident. In Problem 1, if there are K independent sources, each with N states, the number of states for the complete system is N^K . A simplification results if all sources are identical and indistinguishable. In this case the number of states can be reduced to the number of ways of partitioning K things into N groups, which is $(K + N - 1)!/K!(N - 1)!$. Even in this simpler case, we have 5151 states for $K = 100$, $N = 3$. In Problem 2, each of the two traffic sources might have of the order of 100 states, yielding of the order of 10000 states for the complete system. The large numbers are of course the consequence of the fact that the number of states for a system composed of many independent subsystems is (except for special cases) equal to the product of the numbers of states in each subsystem. The complexity of the computation in these systems grows as the *cube* of the number of states. Thus, these numbers are clearly catastrophic.

In the exposition that follows, we propose a general procedure for the analysis of separable Markov-modulated rate processes. The approach is valid irrespective of the specific structure of the subsystems. Two restrictions are crucial to our objective:

- Independence of the underlying Markov processes.
- Separability of the arrival/departure rates.

An additional reversibility assumption will also be made, which is not crucial but considerably simplifies the computational problems. It is shown that, subject to these assumptions, the core of the analysis problem can be reduced to manageable size by decomposition. In order to reduce the computational problems remaining after decomposition, some simple bounds and approximations are obtained. Using this approach, the complexity of the main computation is of the order of $\sum (N^{(k)})^3$ rather than $(\prod N^{(k)})^3$, where $N^{(k)}$ is the number of states in the k th subsystem. Not only does this vastly reduce computational requirements but it also generally results in increased accuracy at critical points in the analysis. In addition and most importantly, the formalism yields considerable insight into the behavior of these

systems. In particular, it reveals the importance of the concept of reversibility in this application, and it explains the considerable power of the methods proposed in [1], [13].

In Section 2 we present the general mathematical model, derive the equations for the buffer occupancy distributions and give the form of their solutions. The general structure of the model is similar to those used in [10] and [15]. The equilibrium probabilities are expressed in terms of the solutions of a related differential equation. We demonstrate certain fundamental properties of the eigenvalues and eigenfunctions of this equation, which are the basis for the decomposition procedures which follow. Section 3 presents bounds and approximations for the equilibrium probabilities. Section 4 deals with separability of the equations of buffer occupancy distribution, yielding eigenvector solutions in the form of Kronecker (tensor) products, and the analytical procedure is applied to Problem 2 in Section 5. Section 6 presents conclusions and points out some extensions of this work which are currently underway.

2. The general model

Consider the system of Figure 3, where \mathcal{S} is the state of a finite irreducible Markov process, C is the maximum service rate, henceforth called the ‘capacity’, $\gamma(\mathcal{S})$ is an arrival rate modulated by \mathcal{S} , $(C - \nu(\mathcal{S}))$ is a service rate modulated by \mathcal{S} , and X (a non-negative continuous random variable) is the buffer content. (Without loss of generality, C , γ and ν are assumed to be non-negative.) The behavior of $X(t)$ in the infinite buffer case is defined by

$$(5) \quad \frac{dX}{dt} = r(\mathcal{S}) - C, \quad X > 0$$

where

$$r(\mathcal{S}) = \gamma(\mathcal{S}) + \nu(\mathcal{S}).$$

Note that as far as the buffer content is concerned, it makes no difference whether the function $\nu(\mathcal{S})$ is considered to be a portion of capacity removed from the maximum capacity C , or arrival rate added to the rate γ . However, it is important to observe that the waiting time is different in the two cases. (If $\nu = 0$, the waiting time for arrivals to a buffer of length x is x/C . However, if $\nu \neq 0$ this waiting time is a random variable.) In what follows, we refer to $r(\mathcal{S})$ as the ‘net arrival rate’.

Equation (5), together with the properties of the modulating process, defines the evolution in time of the state (\mathcal{S}, X) of a bivariate Markov process, where \mathcal{S} takes on

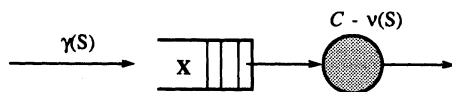


Figure 3. Model of a general Markov-modulated source/server

values s in the set

$$\mathcal{S} = \{s_1, s_2, \dots, s_N\}$$

and X takes on values x in the set $[0, \infty)$ for the infinite buffer case, or $[0, X_{\max}]$ for the finite buffer case. In the latter case, (5) is valid for $0 < X < X_{\max}$.

Let

$$P(t, s, x) = \Pr \{S(t) = s, X(t) \leq x\}.$$

Because the modulating process is finite and irreducible its equilibrium probabilities

$$\pi(s) = \lim_{t \rightarrow \infty} \Pr \{S(t) = s\}$$

exist, are all positive, and are independent of the initial state distribution. They are found from the properties of the modulating process alone.

An important system parameter derived from $\pi(s)$ is the mean value of the net arrival rate

$$\bar{r} = \sum_{s \in \mathcal{S}} \pi(s)r(s)$$

from which we obtain the system traffic intensity

$$\rho = \bar{r}/C.$$

A necessary and sufficient condition for the existence of equilibrium probabilities $F(s, x)$ for the joint process (S, X) in the infinite buffer case is $\rho < 1$. We henceforth assume that this condition is fulfilled, in which case

$$F(s, x) = \lim_{t \rightarrow \infty} P(t, s, x).$$

For a system with infinite buffer

$$(6) \quad \pi(s) = \lim_{x \rightarrow \infty} F(s, x).$$

For a finite buffer system, there will generally be a discontinuity in the function $F(s, x)$ at X_{\max} , so that

$$(7) \quad \pi(s) = P_f(s) + F(s, X_{\max}^-)$$

where

$$F(s, X_{\max}^-) = \lim_{x \rightarrow X_{\max}} F(s, x)$$

and the full buffer equilibrium probabilities $P_f(s)$ are defined as

$$P_f(s) = \lim_{t \rightarrow \infty} \Pr \{S(t) = s, X(t) = X_{\max}\}.$$

As indicated previously, the modulating process is a finite state Markov process. Let $M(s, u)$ be the transition rate from state u to state s for this process, $s \neq u$, and

define

$$M(\mathbf{s}, \mathbf{s}) = - \sum_{\mathbf{u} \neq \mathbf{s}} M(\mathbf{u}, \mathbf{s}).$$

We partition the states into an 'underload' set, \mathcal{S}_u , and an 'overload' set, \mathcal{S}_o , where

$$\mathcal{S}_u = \{s_j : r(s_j) < C\}, \quad \mathcal{S}_o = \{s_j : r(s_j) > C\}.$$

(For reasons stated below, we assume throughout, without loss of generality, that $r(s_j) \neq C$ for all s_j .) In this section we assume further that the states are labelled so that

$$(8) \quad 0 \leq r(s_1) \leq r(s_2) \leq \dots \leq r(s_\omega) < C < r(s_{\omega+1}) \leq \dots \leq r(s_N)$$

where ω and $(N - \omega)$ are respectively the cardinality of the underload and overload sets. This gives

$$\mathcal{S}_u = \{s_1, \dots, s_\omega\}, \quad \mathcal{S}_o = \{s_{\omega+1}, \dots, s_N\}.$$

As can be seen from (1), the buffer occupancy decreases during underload and increases during overload. The Kolmogorov differential equation defining the function $P(t, \mathbf{s}, x)$ for this system is

$$(9) \quad \frac{\partial P}{\partial t} + d(\mathbf{s}) \frac{\partial P}{\partial x} = \sum_{\mathbf{u}} M(\mathbf{s}, \mathbf{u}) P(t, \mathbf{u}, x)$$

where

$$(10) \quad d(\mathbf{s}) = r(\mathbf{s}) - C.$$

To find the equilibrium probabilities we set $\partial P / \partial t = 0$ in (9) to obtain

$$(11) \quad d(\mathbf{s}) F'(\mathbf{s}, x) = \sum_{\mathbf{u}} M(\mathbf{s}, \mathbf{u}) F(\mathbf{u}, x)$$

where $(\cdot)'$ denotes differentiation with respect to x . Equation (11) represents a set of N linear ordinary differential equations which, with suitable boundary conditions, can be solved uniquely for F . Note that if $r(s_j) = C$ for some s_j , the set (11) is singular. In that case, the corresponding equation and the state s_j can be eliminated yielding a lower order system. We avoid this case by imposing the two strict inequalities stated in (8). Note also that if F' is set equal to zero in (11) we obtain N homogeneous linear equations for the equilibrium probabilities $\pi(\mathbf{s})$ of the modulating process.

Denoting the states as integers:

$$s_j = j, \quad j = 1, 2, \dots, N$$

and letting

$$\mathbf{F}(x) = [F(1, x), F(2, x), \dots, F(N, x)]^t$$

$$\mathbf{R} = \text{diag} \{r(j)\}$$

$$\mathbf{D} = \text{diag} \{d(j)\}$$

$$\mathbf{M} = [M(i, j)]$$

where $(\cdot)'$ denotes transpose, (11) becomes

$$(12) \quad DF' = MF$$

where M , an $N \times N$ matrix, is the infinitesimal generator matrix for the underlying Markov process. The solution of (12) takes the form

$$(13) \quad F(x) = \sum_n a_n \varphi_n \exp(z_n x)$$

where

$$\begin{aligned} \varphi_n &= (\varphi_n(s_1), \dots, \varphi_n(s_N))' \\ &= (\varphi_n(1), \dots, \varphi_n(N))' \end{aligned}$$

and each pair $\{z_n, \varphi_n\}$ now satisfies the generalized eigenvalue problem

$$(14) \quad Dz\varphi = M\varphi.$$

(While (14) could be replaced by a standard eigenvalue problem for the matrix $D^{-1}M$, the development that follows is based on the special properties of D and M , which would be obscured in the standard eigenvalue formulation.)

Because M is the generator matrix of an irreducible Markov process, it has one and only one zero eigenvalue, (also a solution of (14)), which will be denoted z_1 . Its eigenvector φ_1 , normalized so that its elements are positive and sum to unity, is

$$\varphi_1 = \pi = (\pi(s_1), \pi(s_2), \dots, \pi(s_N))'$$

the vector of equilibrium probabilities.

Up to this point, we are dealing with a standard problem and solution of (11) subject to given boundary conditions is in theory straightforward. However, it becomes intractable in most cases of practical interest simply because of its size. As mentioned earlier, the number of equations can easily range from hundreds to tens of thousands. In such cases, brute force techniques are generally useless. Computing time, memory requirements numerical accuracy, and even the specification of the equations themselves can all be serious problems. While large computers may make speed and memory requirements minor considerations, inherent problems associated with ill-conditioning and numerical instability will generally cause conventional numerical approaches to break down. The techniques which follow are designed to mitigate problems associated with large N , by exploiting the special structure of the system equations.

2.1. *Reversible processes: properties of the eigenvalues.* Henceforth we make the additional assumption that the underlying Markov process obeys the Kolmogorov condition

$$(15) \quad M(s, u)\pi(u) = M(u, s)\pi(s).$$

Equation (15) is a necessary and sufficient condition for reversibility [11]. Although this assumption is not required for separability, it greatly simplifies the subsequent development. A consequence of reversibility [11], which follows directly from condition (15), is that the infinitesimal generator matrix M in (12) is symmetrized by

the similarity transformation,

$$(16) \quad \tilde{M} = E^{-1}ME$$

where

$$E = \text{diag} \{e(s_j)\}: \quad j = 1, 2, \dots, N$$

and

$$e(s_j) = \sqrt{\pi(s_j)}.$$

The matrix $-\tilde{M}$ is symmetric positive semi-definite. From (16) we see that the change of variable

$$F = E\tilde{F}$$

reduces (12) to the symmetric form

$$D\tilde{F}' = \tilde{M}\tilde{F}$$

and (14) to the symmetric form

$$(17) \quad Dz\tilde{\varphi} = \tilde{M}\tilde{\varphi}$$

where

$$\varphi = E\tilde{\varphi}.$$

We refer to $\tilde{\varphi}$ as a symmetrized eigenvector.

If φ and ψ are real N -vectors we denote their inner product as

$$(\varphi, \psi) = (\psi, \varphi) = \varphi^t \psi.$$

One of the useful consequences of symmetry in (17) is that the eigenvalues z_i are real and a set of N real linearly independent symmetrized eigenvectors $\tilde{\varphi}_n$ can be found which obey the following generalized orthogonality relations:

$$(18) \quad (\tilde{\varphi}_i, D\tilde{\varphi}_j) = 0, \quad i \neq j$$

$$(19) \quad (\tilde{\varphi}_1, D\tilde{\varphi}_1) = (\bar{r} - C)(\tilde{\varphi}_1, \tilde{\varphi}_1)$$

$$(20) \quad (\tilde{\varphi}_i, D\tilde{\varphi}_i) = (1/z_i)(\tilde{\varphi}_i, \tilde{M}\tilde{\varphi}_i), \quad i = 2, 3, \dots, N$$

(see Section A.1 for proof).

Note that $(\tilde{\varphi}_i, D\tilde{\varphi}_i)$ is negative for $z_i > 0$ and positive for $z_i < 0$. Henceforth we assume that the $\{\tilde{\varphi}\}$ are normalized to unity, i.e.

$$(\tilde{\varphi}_i, \tilde{\varphi}_i) = 1, \quad i = 1, 2, \dots, N.$$

For $\tilde{\varphi}_1$, this normalization gives

$$(\tilde{\varphi}_1, D\tilde{\varphi}_1) = \bar{r} - C = C(\rho - 1)$$

and, assuming the appropriate sign for the normalization constant

$$\varphi_1 = E\tilde{\varphi}_1 = \pi.$$

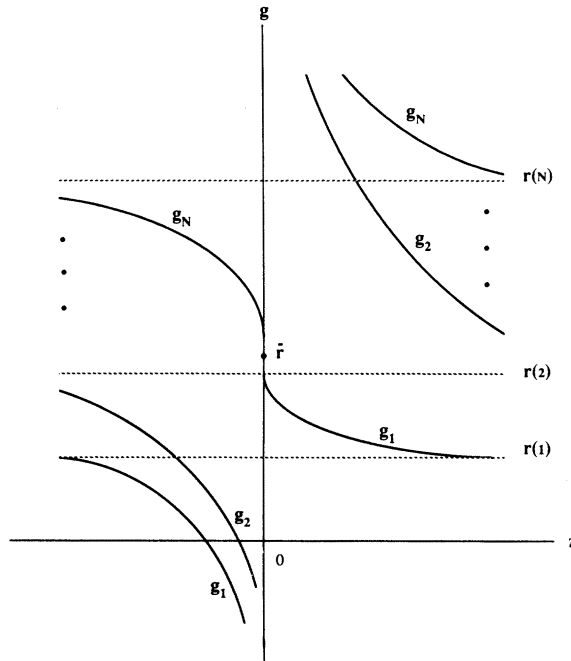


Figure 4. The functions $g_n(z)$

The decomposition techniques described in the subsequent sections are based on certain key relations among the $\{z_n\}$, the $\{\varphi_n\}$ and the system parameters, particularly the capacity C . These relations, proved in Section A.2, are as follows.

For any $C > \bar{r}$ there exists a set of N real solutions, $\{z_n, \varphi_n\}$, of (14), where $z_1 = 0$, z_n is the unique solution of the equation

$$(21) \quad C = g_n(z), \quad n = 2, 3, \dots, N$$

and the set $\{\varphi_n\}$ is linearly independent. The function $g_n(z)$ is the n th eigenvalue of the matrix

$$(22) \quad A(z) = R - M/z, \quad z \neq 0.$$

Letting $\varphi(n; z)$ denote an eigenvector corresponding to $g_n(z)$, the vector $\varphi(n; z_n)$ satisfies (14), so that $\varphi_n = \varphi(n; z_n)$. (The form of the family $\{g_n\}$ is illustrated in Figure 4). Each function $g_n(z)$ is continuous and is continuously differentiable. At each z at which $A(z)$ has distinct eigenvalues its derivative is

$$g_n'(z) = (\bar{\varphi}(n; z) \bar{M} \bar{\varphi}(n; z)) / z^2 < 0, \quad n = 1, 2, \dots, N$$

where

$$\bar{\varphi}(n; z) = E^{-1} \varphi(n; z)$$

and

$$(\bar{\varphi}(n; z), \bar{\varphi}(n; z)) = 1.$$

The functions g_n are indexed so that their asymptotic behavior is as follows:

$$\lim_{z \rightarrow 0^+} g_n(z) = +\infty, \quad n = 2, 3, \dots, N$$

$$\lim_{z \rightarrow 0^-} g_n(z) = -\infty, \quad n = 1, 2, \dots, N-1$$

$$\lim_{z \rightarrow \pm\infty} g_n(z) = r(n), \quad n = 1, 2, \dots, N$$

$$\lim_{z \rightarrow 0^+} g_1(z) = \lim_{z \rightarrow 0^-} g_N(z) = \bar{r}.$$

Note from (8) that the functions $g_n(z)$ are indexed to be increasing functions of n for sufficiently large $|z|$. For future reference this will be termed ‘asymptotically ascending order’. (If $A(z)$ has distinct eigenvalues for all $z \neq 0$ then the ascending order is maintained for all z , as shown in Figure 4.)

The above development is simply an alternative way of determining solutions of (14). Thus, the generalized eigenvalue problem (14) has been replaced by two other problems: a standard eigenvalue problem involving $A(z)$ followed by functional inversion. While this approach may appear convoluted, its usefulness will be apparent in Section 4.

It follows from the above that the $\{z_n\}$ satisfy

$$(23) \quad \begin{aligned} z_n &> 0, & n &= 2, 3, \dots, \omega \\ z_n &< 0, & n &= \omega + 1, \omega + 2, \dots, N. \end{aligned}$$

If $A(z)$ has distinct eigenvalues for all z then the z_n are ordered as follows:

$$(24) \quad z_{\omega+1} \leq \dots \leq z_N < 0 = z_1 < z_2 \leq \dots \leq z_\omega$$

with ω defined in (8). Hence, for each value of C satisfying (8), the general solution of (12) given in (13), is composed of a set of $(N - \omega)$ ‘stable’ modes (those terms in (13) with negative z_n), $(\omega - 1)$ ‘unstable’ modes (those with positive z_n) and the ‘equilibrium’ mode, φ_1 , independent of x .

2.2. Boundary conditions. By using known properties of the function $F(x)$ at the boundaries $x = 0$ and X_{\max} (or $+\infty$), it is possible to determine the coefficients a_n in (13). To this end it is convenient to write (13) in the more compact form

$$(25) \quad F(x) = \Phi e^{Zx} a$$

where

$$\begin{aligned} \Phi &= (\varphi_1, \varphi_2, \dots, \varphi_N) = (\Phi_u, \Phi_s) \\ \Phi_u &= (\varphi_1, \varphi_2, \dots, \varphi_\omega) \\ \Phi_s &= (\varphi_{\omega+1}, \varphi_{\omega+2}, \dots, \varphi_N) \\ e^{Zx} &= \text{diag} \{ \exp(z_n x) : n = 1, 2, \dots, N \} \\ a &= (a_1, a_2, \dots, a_N)^t = (a_u^t, a_s^t)^t. \end{aligned}$$

The partitions Φ_u and Φ_s represent the eigenvectors associated with unstable and stable modes respectively, where the equilibrium mode $\varphi_1 = \pi$ is included with the unstable modes. The vector a is partitioned similarly. Throughout this section the subscripts u (unstable) and s (stable) will be used to denote partitions of a matrix consisting of the first ω columns and last $(N - \omega)$ columns respectively. Similarly, the subscripts u (underload) and o (overload) will be used to denote partitions of a matrix consisting of the first ω rows and last $(N - \omega)$ rows respectively.

Let

$$(26) \quad f = F(0) = \Phi a$$

be the equilibrium empty buffer probability vector, where

$$f = (f(s_1), f(s_2), \dots, f(s_N))^t$$

and

$$f(s) = F(s, 0).$$

We partition f as follows:

$$f = (f_u^t, f_o^t)^t.$$

The partitions f_u and f_o correspond to the underload and overload states respectively. We note that the buffer is never empty while in overload so that

$$(27) \quad f_o = \mathbf{0}.$$

In the infinite buffer case, because $F(x)$ is bounded, the coefficients of the unstable modes must vanish (with the exception of $a_1 = 1$), i.e.,

$$(28) \quad a_u = (1, 0, 0, \dots, 0)^t.$$

Substituting this condition into (26) yields

$$f = \pi + \Phi_s a_s.$$

Letting

$$\pi = (\pi_u^t, \pi_o^t)^t$$

$$\Phi_s = (\Phi_{us}^t, \Phi_{os}^t)^t$$

we have

$$(29) \quad \Phi_{os} a_s = -\pi_o$$

which is a set of $(N - \omega)$ linear equations determining the unknown coefficients of the stable modes, a_s . This approach to finding the coefficients is effective when the number of overload states is relatively small. In the opposite case, where ω is relatively small, the following alternative approach is more suitable. We write (26) in the form

$$(30) \quad a = \Phi^{-1} f$$

where the inverse of Φ can be expressed explicitly using the orthogonality relations (20). From (26) we have

$$(31) \quad a_i = (\tilde{\varphi}_i, DE^{-1}f) / (\tilde{\varphi}_i, D\tilde{\varphi}_i)$$

or

$$(32) \quad (\Phi^{-1})_{ij} = \tilde{\varphi}_j(s_i)d(s_j)/e(s_j)(\tilde{\varphi}_i, D\tilde{\varphi}_i).$$

The last $(N - \omega)$ rows of (30) give the stable mode coefficients, a_s , in terms of the empty buffer probabilities f , of which those for the underload states, f_u , remain to be determined. These are found using the first ω rows of (30):

$$(33) \quad (\Phi^{-1})_{uu}f_u = a_u$$

where a_u is given in (28). Equation (33) represents a set of ω linear equations for the empty buffer probabilities, f_u , the solution of which must be substituted into (30) to find a_s .

Since solution of L linear equations requires of the order of L^3 operations, the computational complexity of the above procedures is dominated by the solution of (29) in the former case, or (33) in the latter case.

In the finite buffer case, all coefficients a_i are generally present in the solution (13). Furthermore, the coefficient of the constant term in (13) is no longer unity. Instead, we have (7) relating $\pi(s)$ to the full buffer probability $P_f(s)$.

Let

$$f_{\max} = (F(s_1, X_{\max-}), F(s_2, X_{\max-}), \dots, F(s_N, X_{\max-}))^t$$

$$P_f = (P_f(s_1), P_f(s_2), \dots, P_f(s_N))^t.$$

Then in addition to the boundary condition (26), which holds for all systems, we have from (26) and (7)

$$(34) \quad f_{\max} = \Phi \exp(ZX_{\max})a = \pi - P_f.$$

Noting that the buffer can never be full while in an underload state, we have

$$P_f(s) = 0, \quad s \in \mathcal{S}_u.$$

These additional boundary conditions lead to $2N$ linear equations to be solved for a .

In Section 3 we derive some approximations based on these relations, which greatly reduce the computational burden.

3. Coefficient bounds, approximations

In most problems involving large systems a determination of the exact probability distribution $F(x)$ is unnecessary and/or computationally infeasible. In such cases it is useful to obtain simple approximations using only a few terms in the summation (13). We consider this problem now. Our results are illustrated with a numerical

example in Section 5. The performance parameters of interest in systems of the type described in Section 1 are typically scalar functions of $F(\mathbf{s}, x)$. Two functions which are particularly important are

$$(35) \quad G(x) = \Pr \{X > x\} = (\mathbf{1}, \boldsymbol{\pi} - \mathbf{F})$$

where

$$\mathbf{1} = (1, 1, \dots, 1)^t$$

and

$$(36) \quad H(x) = (\mathbf{1}, \mathbf{R}(\boldsymbol{\pi} - \mathbf{F}))/\bar{r}.$$

In most well-designed systems the size X_{\max} of the (finite) buffer is chosen so that $(\mathbf{1}, \mathbf{P}_f)$, the probability of a full buffer, is small. In such cases $G(X_{\max}) \approx (\mathbf{1}, \mathbf{P}_f)$, and $G(X_{\max})$ overbounds $(\mathbf{1}, \mathbf{P}_f)$, so that the infinite buffer model can be used to determine the probability of buffer overflow in these cases. The function $H(x)$ defined in (36) is the fraction of flow that arrives to a buffer whose content exceeds x , or equivalently, the probability that an arriving data unit sees a buffer of content exceeding x . Thus, $H(X_{\max})$ approximates the fraction of arriving flow that is 'blocked' by a finite buffer. For a system with constant service rate ($\nu = 0$), $H(CT)$ represents the fraction of flow that is delayed a time exceeding T , an important performance parameter. We seek approximations to G and H which can be obtained with a minimal amount of computation.

Consider first $G(x)$. In view of the above remarks we limit the discussion to the infinite buffer case and assume that the $\{z_n\}$ are ordered as in (24) with repeated eigenvalues permitted. Equation (13) can then be rewritten in the form

$$(37) \quad \mathbf{F}(x) = \sum_{n=\omega+1}^N a_n \boldsymbol{\varphi}_n \exp(z_n x) + \boldsymbol{\pi}.$$

Substituting (35) into (37) we have

$$(38) \quad G(x) = \sum_{n=\omega+1}^N b_n \exp(z_n x)$$

where

$$(39) \quad b_n = -a_n(\mathbf{1}, \boldsymbol{\varphi}_n).$$

If the number of overload states is large, the computation of (38) may be both onerous and unnecessary. Instead, a good approximation of $G(x)$ may usually be obtained by truncating the series to a small number of terms. In the context of well-designed systems, mentioned above, it is the larger values of x that are of interest. For these values the dominating terms are those for large n , since the modes with lower indices generally decay rapidly with x . The truncated series retaining m terms of highest index will be denoted

$$(40) \quad G_m(x) = \sum_{n=N-m+1}^N b_n \exp(z_n x), \quad 1 \leq m \leq N - \omega.$$

The exponents z_n in (40) are solutions of (14), which are readily computed for separable systems using the decomposition methods of Section 4. In some cases they can be expressed in closed form. The exponent z_N , in the single exponential approximation ($m = 1$), is the most important parameter in characterizing system performance, since it indicates the asymptotic behavior of all performance measures of interest for large x . (The function relating z_N to ρ in matrix-geometric systems, termed the 'caudal characteristic', is studied by Neuts [19].) Of course, if behavior for small values of x is of interest then the rapidly decaying modes are important [24], [17]. The coefficients b_n in (40) are found by substituting the expression for a_n in (31) into (39) to give

$$(41) \quad \begin{aligned} b_n &= (\mathbf{1}, \varphi_n)(\tilde{\varphi}_n, -D\tilde{f})/(\tilde{\varphi}_n, D\tilde{\varphi}_n) \\ &= (\mathbf{1}, \varphi_n)(\tilde{\varphi}_n, -D\tilde{f})z_n/(\tilde{\varphi}_n, \tilde{M}\tilde{\varphi}_n) \end{aligned}$$

where

$$\tilde{f} = E^{-1}f.$$

In what follows we assume that the sign of the normalization constant of $\tilde{\varphi}_n$ is chosen so that $(\mathbf{1}, \varphi_n) \geq 0$, which implies that $(\mathbf{1}, \varphi_n)/(\tilde{\varphi}_n, D\tilde{\varphi}_n) \geq 0$. (Our previous assumption that $\tilde{\varphi}_n$ is normalized to unity still leaves the sign undetermined.) Because of the dependence of b_n on the empty buffer probability vector f , which is unknown, the coefficients b_n cannot be computed on a term by term basis as is the case for the exponents z_n . If it happens that the methods of Section 2.2.1 for obtaining a_n (either directly, or indirectly via f) are unsuitable because of the size of the system, it is possible to obtain some good upper and lower bounds on b_n by using known constraints on f . Recall from Section 2 that

$$(42) \quad 0 \leq f(s) \leq \pi(s), \quad s \in \mathcal{S}_u,$$

$$(43) \quad f(s) = 0, \quad s \in \mathcal{S}_o.$$

Also, it is easily shown (see Section A.3) that

$$(44) \quad \sum_s -d(s)f(s) = C(1 - \rho).$$

Now let

$$(45) \quad \begin{aligned} \bar{b}_n &= \frac{(\mathbf{1}, \varphi_n)}{(\tilde{\varphi}_n, D\tilde{\varphi}_n)} \max_f (\tilde{\varphi}_n, -D\tilde{f}) \\ &= \frac{(\mathbf{1}, \varphi_n)}{(\tilde{\varphi}_n, D\tilde{\varphi}_n)} \max_f \sum_{s \in \mathcal{S}_u} \frac{-\varphi_n(s)d(s)f(s)}{\pi(s)} \end{aligned}$$

$$(46) \quad \begin{aligned} b_n &= \frac{(\mathbf{1}, \varphi_n)}{(\tilde{\varphi}_n, D\tilde{\varphi}_n)} \min_f (\tilde{\varphi}_n, -D\tilde{f}) \\ &= (\mathbf{1}, \varphi_n) \min_f \sum_{s \in \mathcal{S}_u} \frac{-\varphi_n(s)d(s)f(s)}{\pi(s)} \end{aligned}$$

where the maximization (minimization) is effected subject to the constraints (42), (43), (44). The coefficients \bar{b}_n and \underline{b}_n are respectively upper and lower bounds on b_n . To evaluate \bar{b}_n , the terms in the summation in (45) are arranged so that the quantities $\varphi_n(s)/\pi(s)$ appear in descending order, and $f(s)$ is assigned its maximum possible value (which is $\leq \pi(s)$) in each term until the condition (44) is satisfied. Similarly, the minimum \underline{b}_n , is attained when the maximum possible value of $f(s)$ is assigned to each term in (46) in *ascending* order of the quantities $\varphi_n(s)/\pi(s)$ until condition (44) is satisfied. A slightly looser upper bound \hat{b}_n , on b_n can be found more simply without reordering the summation:

$$(47) \quad \hat{b}_n = \frac{(\mathbf{1}, \varphi_n)}{(\bar{\varphi}_n, D\bar{\varphi}_n)} \sum_{s \in \mathcal{S}_u} -\max [0, \varphi_n(s)] d(s).$$

In (47) the upper bound is found by assigning $f(s)$ its maximum value in each term while ignoring the equality constraint (44). Good approximations to $G(x)$, which will overbound and underbound it for sufficiently large x , are found by replacing the coefficients b_n in (38) by \bar{b}_n or \underline{b}_n respectively. Thus, define

$$(48) \quad \bar{G}_m(x) = \sum_{n=N-m+1}^N \bar{b}_n \exp(z_n x),$$

$$(49) \quad \underline{G}_m(x) = \sum_{n=N-m+1}^N \underline{b}_n \exp(z_n x).$$

The advantage of using this technique is that the computation can proceed one term at a time, until a satisfactory approximation is obtained. It also explicitly exhibits the exponents, which dominate in determining system performance.

Approximations for $H(x)$ are obtained in a manner exactly paralleling those for $G(x)$ by replacing each occurrence of the quantity $(\mathbf{1}, \varphi_n)$ in the above development by $(\mathbf{1}, R\varphi_n)/\bar{r}$. Note that this requires replacing the condition, $(\mathbf{1}, \varphi_n) \geq 0$, by $(\mathbf{1}, R\varphi_n) \geq 0$, which may require reversing the signs of some of the normalization constants.

If $G(x)$ and $H(x)$ are to be evaluated at given set of values of x , the above approximation can be refined [23].

4. Separability and decomposition

As mentioned earlier, a direct application of the formalism presented in Section 2 leads to formidable computational problems when the system size is very large. However, large systems are very often composed of combinations of independent and 'separable' smaller systems. We now show that in this case the bulk of the computation can be performed at the level of each subsystem. This can produce orders of magnitude reduction in the computational burden and increased accuracy in the numerical algorithms.

Let the system (14) be composed of K independent subsystems, where $s^{(k)}$ denotes the state of the k th subsystem, and the K -tuple

$$\mathbf{s} = (s^{(1)}, s^{(2)}, \dots, s^{(K)})$$

denotes the combined or 'global' system state. The k th subsystem ($k = 1, \dots, K$) is assumed to be a $N^{(k)}$ -state, irreducible, reversible Markov process whose states will be denoted by the integers

$$s^{(k)} \in \{1, 2, \dots, N^{(k)}\} = \mathcal{S}^{(k)}$$

where

$$N = \prod_k N^{(k)}$$

so that each state \mathbf{s}_j is an integer K -tuple. The indices will be ordered lexicographically, i.e.

$$\begin{aligned} \mathbf{s}_1 &= (1, 1, \dots, 1) \\ \mathbf{s}_2 &= (1, 1, \dots, 2) \\ &\vdots \\ \mathbf{s}_N &= (N^{(1)}, N^{(2)}, \dots, N^{(K)}). \end{aligned}$$

Let $\mathbf{M}^{(k)}$ denote the matrix of transition rates for the k th subsystem. The buffer content X is now the accumulation of arrivals from the K subsystems, and the net arrival rate $r(\mathbf{s})$ in (1) is composed of the superposition of net arrival rates $r^{(k)}$ for each subsystem:

$$(50) \quad r(\mathbf{s}) = \sum_k r^{(k)}(s^{(k)})$$

where

$$(51) \quad r^{(k)}(s^{(k)}) = \gamma^{(k)}(s^{(k)}) + \nu^{(k)}(s^{(k)})$$

and we define a matrix of net arrival rates for each subsystem as

$$(52) \quad \mathbf{R}^{(k)} = \text{diag} \{r^{(k)}(1), r^{(k)}(2), \dots, r^{(k)}(N^{(k)})\}.$$

It is assumed that the states of each subsystem are labelled in such a way that

$$r^{(k)}(1) \leq r^{(k)}(2) \leq \dots \leq r^{(k)}(N^{(k)}), \quad k = 1, 2, \dots, K.$$

The lexicographic ordering of the global states will not necessarily label them so that condition (8) is satisfied. Nevertheless, all of the development of the previous sections still holds, subject to a minor modification of notation in Section 2.2. There, in order to maintain the proper row partitioning of the various matrices, corresponding to underload and overload states, it may be necessary to relabel some of the global states.

Independence of the subsystems implies that the allowable state transitions correspond to changes in only one element of the state K -tuple at a time, with the state transition rates for the k th subsystem defined by the matrix $M^{(k)}$. We will refer to condition (50) together with the independence assumption as ‘separability conditions’. Subject to the above separability conditions it is easily verified that the matrices D and M in the global equilibrium equation (12) and the related eigenvalue equation (14) are given by

$$(53) \quad M = M^{(1)} \oplus M^{(2)} \oplus \dots \oplus M^{(K)}$$

and

$$(54) \quad D = R - CI$$

where

$$(55) \quad R = R^{(1)} \oplus R^{(2)} \oplus \dots \oplus R^{(K)}.$$

In the above equations and those that follow \oplus and \otimes denote the Kronecker sum and product respectively (see Section A.4). Also, I denotes an identity matrix of appropriate size.

The Kronecker sum form of (53), and (55) suggests a decomposition method for obtaining solutions of (14) by working at the level of the subsystems. The method includes as special cases results appearing in [13] and [15].

Paralleling the development of Section 2, we assume that each subsystem is an irreducible reversible Markov process. For each k , the zero eigenvalue of $M^{(k)}$, will be denoted z_1 . Its eigenvector $\varphi_1^{(k)}$, (assumed to be normalized so that its elements are positive and sum to unity), is

$$\pi^{(k)} = (\pi^{(k)}(1), \pi^{(k)}(2), \dots, \pi^{(k)}(N^{(k)}))'$$

the vector of equilibrium probabilities. The equilibrium vector π for the global process M is then given by

$$\pi = \pi^{(1)} \otimes \pi^{(2)} \otimes \dots \otimes \pi^{(K)}.$$

In terms of these equilibrium probabilities, we define the mean net arrival rate for the k th subsystem as

$$\bar{r}^{(k)} = \sum_j \pi^{(k)}(j) r^{(k)}(j),$$

so that for the global system

$$\bar{r} = \sum_k \bar{r}^{(k)}.$$

Letting

$$e^{(k)}(j) = \sqrt{\pi^{(k)}(j)}, \quad j = 1, 2, \dots, N^{(k)}$$

$$E^{(k)} = \text{diag} \{e^{(k)}(1), \dots, e^{(k)}(N^{(k)})\}$$

we may symmetrize the matrix $\mathbf{M}^{(k)}$ by the transformation

$$(56) \quad \tilde{\mathbf{M}}^{(k)} = (\mathbf{E}^{(k)})^{-1} \mathbf{M}^{(k)} \mathbf{E}^{(k)}.$$

We are now prepared to state the relations among the $\{z_n\}$, the $\{\varphi_n\}$ and the system parameters, in terms of the characteristics of the subsystems. (See Section A.5 for proof.)

For any $C > \bar{r}$ there exists a set of N real solutions, $\{z_n, \varphi_n\}$, of (14), where $z_1 = 0$, z_n is the unique solution of the equation

$$(57) \quad C = g_n(z), \quad n = 2, 3, \dots, N$$

and the set $\{\varphi_n\}$ is linearly independent. The functions $g_n(z)$ are defined as

$$(58) \quad g_n(z) = \sum_{k=1}^K g^{(k)}(i_n^{(k)}; z), \quad n = 1, 2, \dots, N$$

where $g^{(k)}(i; z)$ is the i th eigenvalue of the matrix,

$$(59) \quad \mathbf{A}^{(k)}(z) = \mathbf{R}^{(k)} - \mathbf{M}^{(k)}/z, \quad z \neq 0$$

indexed in asymptotically ascending order, and $\varphi^{(k)}(i; z)$ is a corresponding eigenvector. Letting $\varphi_n^{(k)}$ denote $\varphi^{(k)}(i_n^{(k)}; z_n)$, the global eigenvectors are expressed as

$$(60) \quad \varphi_n = \varphi_n^{(1)} \otimes \varphi_n^{(2)} \otimes \dots \otimes \varphi_n^{(K)}.$$

The g_n are indexed in lexicographic order of the arguments $i_n^{(k)}$, i.e.

$$\begin{aligned} g_1(z) &= g^{(1)}(1; z) + g^{(2)}(1; z) + \dots + g^{(K)}(1; z) \\ g_2(z) &= g^{(1)}(1; z) + g^{(2)}(1; z) + \dots + g^{(K)}(2; z) \\ &\vdots \\ g_N(z) &= g^{(1)}(N^{(1)}; z) + g^{(2)}(N^{(2)}; z) + \dots + g^{(K)}(N^{(k)}; z). \end{aligned}$$

Each function $g^{(k)}(i; z)$ is continuous and is continuously differentiable. At each z at which $\mathbf{A}^{(k)}(z)$ has distinct eigenvalues its derivative is

$$(61) \quad g^{(k)}(i; z)' = (\tilde{\varphi}^{(k)}(i; z), \tilde{\mathbf{M}}^{(k)} \tilde{\varphi}^{(k)}(i; z))/z^2, \quad i = 1, 2, \dots, N^{(k)}$$

where

$$\tilde{\varphi}^{(k)}(i; z) = (\mathbf{E}^{(k)})^{-1} \varphi^{(k)}(i; z).$$

The asymptotic behavior of $g^{(k)}$ is as follows:

$$(62) \quad \lim_{z \rightarrow 0^+} g^{(k)}(i; z) = +\infty, \quad i = 2, 3, \dots, N^{(k)},$$

$$(63) \quad \lim_{z \rightarrow 0^-} g^{(k)}(i; z) = -\infty, \quad i = 1, 2, \dots, N^{(k)} - 1,$$

$$(64) \quad \lim_{z \rightarrow \pm\infty} g^{(k)}(i; z) = r^{(k)}(i), \quad i = 1, 2, \dots, N^{(k)},$$

$$(65) \quad \lim_{z \rightarrow 0^+} g^{(k)}(1; z) = \lim_{z \rightarrow 0^-} g^{(k)}(N^{(k)}; z) = \bar{r}^{(k)}.$$

In the above development, each choice of index, $n > 1$, corresponds to the selection of a unique set of terms $g^{(k)}(i_n^{(k)}; z)$ on the right-hand side of (58), whose sum is the function $g_n(z)$, which must be inverted to yield z_n . With z_1 defined to be zero, the remaining $N - 1$ combinations of terms are accounted for as n runs from 2 to N .

The sign of the $\{z_n\}$ can be predicted in a manner analogous to that used when the system is analyzed in its global form. Identifying the K -tuple of indices $i_n^{(k)}$ with the integer K -tuple representation of a global state s_n :

$$s_n = (i_n^{(1)}, i_n^{(2)}, \dots, i_n^{(K)})$$

it follows from the asymptotic behavior of $g_n(z)$ that the non-zero $\{z_n\}$ satisfy

$$(66) \quad \begin{aligned} z_n > 0, & \quad s_n \in \mathcal{S}_u \\ z_n < 0, & \quad s_n \in \mathcal{S}_o. \end{aligned}$$

However, an ordering of the type of (24) cannot be predicted in terms of the properties of the subsystems when decomposition is used.

To make the relation between the global system and the subsystems more explicit, note that $g^{(k)}(i_n^{(k)}; z_n)$ is an eigenvalue and $\boldsymbol{\varphi}_n^{(k)}$ an associated eigenvector for $\mathbf{A}^{(k)}(z_n)$, so that

$$(67) \quad g^{(k)}(i_n^{(k)}; z_n) \boldsymbol{\varphi}_n^{(k)} = \mathbf{A}^{(k)}(z_n) \boldsymbol{\varphi}_n^{(k)}.$$

Using (59) to rearrange terms in (67) we have

$$(68) \quad \mathbf{D}_n^{(k)} z_n \boldsymbol{\varphi}_n^{(k)} = \mathbf{M}^{(k)} \boldsymbol{\varphi}_n^{(k)}, \quad k = 1, 2, \dots, K$$

where

$$(69) \quad \mathbf{D}_n^{(k)} = \mathbf{R}^{(k)} - g^{(k)}(i_n^{(k)}; z_n) \mathbf{I}$$

and

$$\mathbf{D} = \mathbf{D}_n^{(1)} \oplus \mathbf{D}_n^{(2)} \oplus \dots \oplus \mathbf{D}_n^{(K)}, \quad n = 2, 3, \dots, N.$$

Note from (69) that $g^{(k)}(i_n^{(k)}; z_n)$ plays the role of a (fictitious) channel capacity in the k th subsystem.

It can be seen from (68) that $\{z_n, \boldsymbol{\varphi}_n^{(k)}\}$ satisfies a generalized eigenvalue problem analogous to (14) at the level of the k th subsystem, so that $\boldsymbol{\varphi}_n^{(k)}$ is in fact a 'subsystem eigenvector'. The value z_n , however, is a *global* eigenvalue, common to all subsystems as well as the global system. Certain relations between global and subsystem inner products are useful for reducing the computational burden. Assuming that the symmetrized subsystem eigenvectors,

$$\tilde{\boldsymbol{\varphi}}_n^{(k)} = (\mathbf{E}^{(k)})^{-1} \boldsymbol{\varphi}_n^{(k)}$$

are normalized to unity:

$$(\tilde{\boldsymbol{\varphi}}_i^{(k)}, \tilde{\boldsymbol{\varphi}}_i^{(k)}) = 1, \quad i = 1, 2, \dots, N^{(k)}$$

they obey the following relations:

$$(70) \quad \boldsymbol{\varphi}_1^{(k)} = \mathbf{E}^{(k)} \tilde{\boldsymbol{\varphi}}_1^{(k)} = \boldsymbol{\pi}^{(k)}$$

$$(71) \quad (\tilde{\boldsymbol{\varphi}}_i^{(k)}, \mathbf{D}_i^{(k)} \boldsymbol{\varphi}_i^{(k)}) = z_i^{-1} (\tilde{\boldsymbol{\varphi}}_i^{(k)}, \tilde{\mathbf{M}}^{(k)} \tilde{\boldsymbol{\varphi}}_i^{(k)}), \quad i = 2, 3, \dots, N^{(k)}.$$

Furthermore, it follows from the properties of Kronecker sums and products (Section A.4) that the global inner product [20] can be decomposed into a sum of subsystem inner products:

$$(72) \quad \begin{aligned} (\tilde{\boldsymbol{\varphi}}_i, \mathbf{D} \tilde{\boldsymbol{\varphi}}_i) &= \sum_k (\tilde{\boldsymbol{\varphi}}_i^{(k)}, \mathbf{D}_i^{(k)} \tilde{\boldsymbol{\varphi}}_i^{(k)}) \\ &= \sum_k (\tilde{\boldsymbol{\varphi}}_i^{(k)}, \mathbf{R}_i^{(k)} \tilde{\boldsymbol{\varphi}}_i^{(k)}) - C, \quad i = 2, 3, \dots, N. \end{aligned}$$

To clarify the decomposition procedure and to compare its complexity to the 'brute force' technique of analyzing the global system directly, let us examine the computational steps in more detail. We note first that a solution of (57) for z cannot generally be found in closed form. Thus some sort of approximation procedure is required. The fact that the function g_n is well-behaved, with its derivative directly available, suggests that a second-order technique such as Newton's method is appropriate. Letting $z_n(i)$ be the i th iterate for the solution z_n of (57), Newton's method is defined as

$$(73) \quad z_n(i+1) = z_n(i) - \frac{g_n(z_n(i)) - C}{g_n'(z_n(i))}.$$

Each term $g^{(k)}(i_n^{(k)}; z_n(i))$ in $g_n(z_n(i))$ is an eigenvalue of the matrix $\mathbf{A}^{(k)}(z_n(i))$ of the appropriate order, as specified above. The derivative g_n' is evaluated as a sum of terms $g^{(k)}(i_n^{(k)}; z_n(i))'$, where each term is defined in (61). Evaluation of the right-hand side of (73) for a given index n and iterate $z_n(i)$ thus requires the solution of K (small) eigenvalue problems, one for each subsystem. A single pair, $\{g^{(k)}, \tilde{\boldsymbol{\varphi}}^{(k)}\}$, is selected from each subsystem to evaluate g_n and g_n' . At convergence, $z_n(i) \approx z_n$ and $\tilde{\boldsymbol{\varphi}}^{(k)}(i_n^{(k)}; z_n(i)) \approx \tilde{\boldsymbol{\varphi}}_n^{(k)}$, which are the solutions of (68). The eigenvector $\boldsymbol{\varphi}_n$ for the global system is found as the Kronecker product of the subsystem eigenvectors $\boldsymbol{\varphi}_n^{(k)}$ as prescribed in (60). To avoid the singularity of g_n at the origin, the initial iterate must be chosen of the correct sign, predicted by (66). Newton's method exhibits quadratic convergence, so that very few iterations are required given a reasonable initial trial. (In some typical numerical experiments 1–4 iterations were required for convergence.)

To compare computational complexity we consider the infinite buffer case only, and recall that $N - \omega$ modes must be calculated for an exact solution. However, as pointed out in Section 3, it is normally sufficient to compute only a few terms in the summation for $\mathbf{F}(x)$. In this case a fair comparison of the brute force method versus decomposition can be based on the comparative complexity of computing a single mode in (13). (This excludes from consideration the coefficients a_n , which require

very little computation if they are handled using the bounds of Section 3.) The complexity of the eigenvalue problem for the k th subsystem is of the order of $(N^{(k)})^3$, so that complexity of computation for one mode using decomposition is of the order of $\sum (N^{(k)})^3$. This is to be compared to a complexity of the order of $\prod (N^{(k)})^3$ by brute force. As an example, taking $K=3$ and $N^{(k)}=10$ for each subsystem, the comparison is 3×10^3 versus 10^9 . Decomposition therefore results in dramatic savings, especially when N is large. It should be noted, however, that the brute force computation gives *all* modes with a complexity of $\prod (N^{(k)})^3$, while decomposition requires of the order of $N \sum (N^{(k)})^3$ operations to obtain all N modes. Thus, if all modes are to be computed the comparison is 3×10^6 versus 10^9 , still very favorable to decomposition.

5. An application

The analytical and computational procedures discussed in Sections 2, 3 and 4 are best illustrated through an example. In this section we apply the decomposition technique to the model of Problem 2 (Section 1.2) and obtain various performance curves. However, our main objective here is to demonstrate the performance of the *algorithms* rather than the performance of the system. In computing the performance curves we compare exact methods to approximation and bounding techniques, and thereby obtain information on the tightness of the bounds.

The buffer dynamics for the model of Problem 2 are defined in (2). To proceed with the analysis it is necessary to specify the statistical characteristics of the traffic sources. Recall that there are two traffic classes: class x , which is buffered, and class y , which is blocked. It will be assumed that each class is made up of a collection of independent sources which alternate randomly between active and inactive periods, where the activity process for each source is generated by a two-state Markov process.

Let

N_x = the number of class x sources

N_y = the number of class y sources

i_x = the number of active class x sources

i_y = the number of active class y sources

γ_x = the traffic generated by an active source in class x (data units/s)

γ_y = the traffic generated by an active source in class y

λ_x = the transition rate from the inactive to the active state for a source in class x
(s^{-1})

λ_y = the transition rate from the inactive to the active state for a source in class y

μ_x = the transition rate from the active to the inactive state for a source in class x

μ_y = the transition rate from the active to the inactive state for a source in class y

This system can be separated into two independent subsystems where, using the

terminology of Section 4, the system state is defined as

$$\mathbf{s} = (s^{(1)}, s^{(2)}).$$

The subsystem states are

$$s^{(1)} = i_x + 1, \quad s^{(2)} = i_y + 1$$

with

$$N^{(1)} = N_x + 1, \quad N^{(2)} = N_y + 1, \quad N = (N_x + 1)(N_y + 1).$$

Equation (2) then translates to the form of (5), where

$$r(\mathbf{s}) = r^{(1)}(s^{(1)}) + r^{(2)}(s^{(2)})$$

and

$$r^{(1)}(s^{(1)}) = \gamma^{(1)}(s^{(1)}) = (s^{(1)} - 1)\gamma_x, \quad r^{(2)}(s^{(2)}) = \min[\alpha C, (s^{(2)} - 1)\gamma_y].$$

Equations (3) and (4) become

$$\begin{aligned} \frac{dL}{dt} &= \max[(s^{(2)} - 1)\gamma_y - \alpha C, 0], & X > 0 \\ &= \max[(s^{(2)} - 1)\gamma_y + (s^{(1)} - 1)\gamma_x - C, 0], & X = 0. \end{aligned}$$

Letting

$$\lambda^{(1)} = \lambda_x, \quad \lambda^{(2)} = \lambda_y, \quad \mu^{(1)} = \mu_x, \quad \mu^{(2)} = \mu_y,$$

the non-zero elements of the transition rate matrices for the two modulating processes (which are tridiagonal) are

$$M^{(k)}(i, i + 1) = \mu^{(k)}i$$

$$M^{(k)}(i + 1, i) = \lambda^{(k)}(N^{(k)} - i), \quad k = 1, 2, \quad \text{and} \quad i = 1, 2, \dots, N^{(k)} - 1$$

$$M^{(k)}(i, i) = \lambda^{(k)}(i - N^{(k)}) + (1 - i)\mu, \quad k = 1, 2, \quad \text{and} \quad i = 1, 2, \dots, N^{(k)}.$$

It is worth noting that the first subsystem (but not the second) is itself separable. It is, in fact, a system of the form analyzed in [1], whose eigenvalues and eigenvectors can be expressed in closed form.

The above parameters completely define the model in a form suitable for application of the decomposition methods of Section 4. A program was written in FORTRAN to execute the various steps in the algorithm for computing the equilibrium probabilities and evaluating system performance. Quantities of interest in this application are the functions G and H (defined in Section 3). Recall that the function $G(x)$ represents the probability that the class x buffer occupancy exceeds a value x and $H(x)$ represents the fraction of class x traffic that arrives to a buffer whose occupancy exceeds x . Also of interest is the fractional loss f_y for the traffic in class y . This quantity is defined as

$$f_y = \frac{E\{dL/dt\}}{N_y \lambda_y \gamma_y / (\lambda_y + \mu_y)}$$

where $E\{\cdot\}$ denotes expectation. The parameter α partitions the channel, reserving

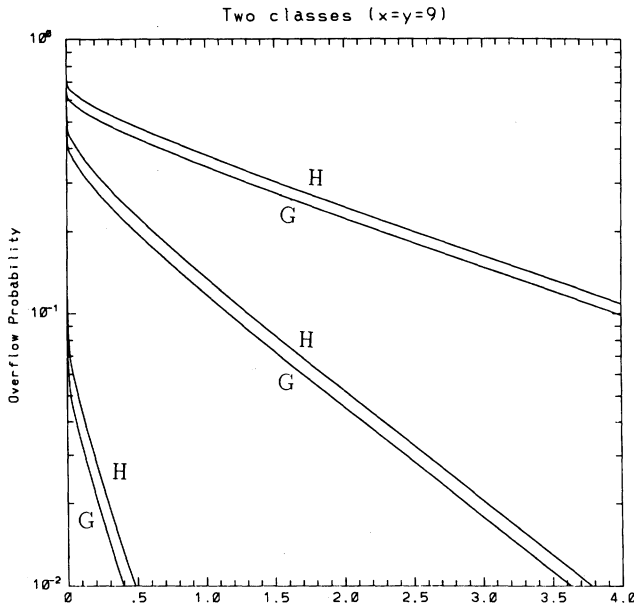


Figure 5. Probabilities of overflow G and H vs. buffer size for $N_x = N_y = 9$. For $\alpha = 1, 0.5,$ and 0.3 , the traffic intensity ρ is set to $0.91, 0.87,$ and 0.74 , respectively

a fraction $(1 - \alpha)$ of the capacity for class x and the rest for class y . Thus it is of interest to explore the tradeoff between delay for class x and blocking for class y as a function of α . This is done through a family of curves in Figures 5 and 6, with system parameters

$$\begin{aligned}
 C &= 9.9 \\
 N^{(1)} &= N^{(2)} = 10 \\
 \lambda^{(1)} &= \lambda^{(2)} = \mu^{(1)} = \mu^{(2)} = \gamma^{(1)} = \gamma^{(2)} = 1.
 \end{aligned}$$

Note that as α increases from 0 (complete priority to class x) to 1 (complete priority to class y), the traffic intensity in the class x buffer increases, thereby increasing the buffer content (Figure 5), but reducing the fraction of blocked class y traffic (Figure 6). Because the higher arrival rates are positively correlated with higher buffer occupancies, the values of H in Figure 5 are consistently larger than G . The quantities $H(CT)$ and $H((1 - \alpha)CT)$ give lower and upper bounds respectively on the probability that buffering delay exceeds T . However, exact calculation of buffering delay requires an additional (and difficult) analytical step, since the capacity available to the class x traffic is a random variable.

The bulk of the computation in this problem is the determination of the equilibrium probabilities. Table 1 shows the essential quantities involved in this computation for the case $\alpha = 0.3$ and $C = 9.9$ with the remaining system parameters as stated above. The (class x) traffic intensity is $\rho = 0.744$, and the fractional loss for

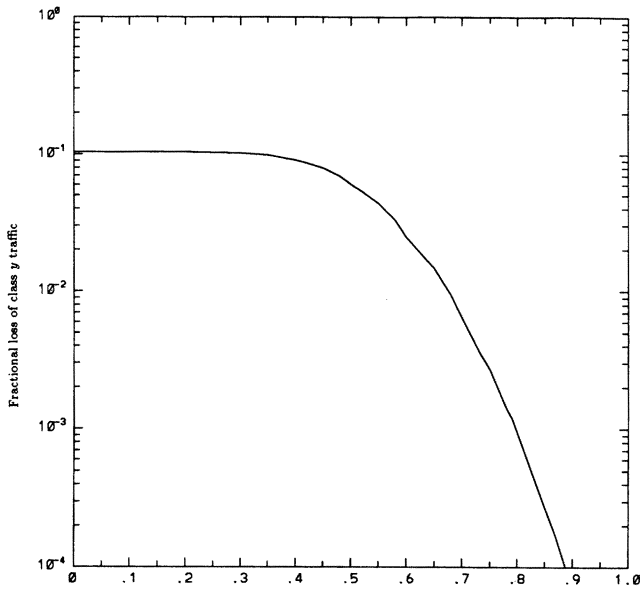


Figure 6. Fractional loss of class y traffic vs α for $N_x = N_y = 9$

TABLE 1

n	z_n	b_n	b_n	\bar{b}_n	\hat{b}_n
100	-0.3167E + 01	0.2990E - 01	0.3125E - 01	0.3273E - 01	0.3273E - 01
99	-0.4776E + 01	0.6575E - 04	0.1545E - 03	0.2547E - 03	0.3981E - 03
98	-0.6338E + 01	0.4122E - 05	0.3696E - 04	0.5291E - 04	0.8742E - 04
92	-0.7801E + 01	0.2128E - 01	0.2294E - 01	0.2440E - 01	0.2440E - 01
97	-0.7893E + 01	0.1574E - 05	0.1463E - 04	0.2145E - 04	0.3191E - 04
96	-0.9436E + 01	0.5655E - 06	0.7285E - 05	0.1073E - 04	0.1564E - 04
91	-0.1046E + 02	0.1190E - 03	0.3324E - 03	0.6176E - 03	0.8334E - 03
95	-0.1095E + 02	0.1210E - 06	0.3995E - 05	0.6267E - 05	0.8163E - 05
94	-0.1240E + 02	-0.5361E - 07	0.1938E - 05	0.3090E - 05	0.4307E - 05
90	-0.1328E + 02	-0.1011E - 04	0.1066E - 03	0.1766E - 03	0.2745E - 03
89	-0.1616E + 02	-0.6972E - 05	0.4864E - 04	0.8426E - 04	0.1194E - 03
93	-0.1624E + 02	0.1866E - 05	0.3162E - 04	0.3907E - 04	0.4303E - 04
88	-0.1906E + 02	-0.4133E - 05	0.2482E - 04	0.4321E - 04	0.6674E - 04
87	-0.2189E + 02	-0.3189E - 05	0.1234E - 04	0.2630E - 04	0.3752E - 04
86	-0.2456E + 02	-0.2589E - 05	0.4521E - 05	0.1069E - 04	0.1884E - 04
83	-0.1310E + 03	0.3903E - 01	0.4529E - 01	0.4969E - 01	0.4969E - 01
82	-0.1713E + 03	0.1592E - 03	0.1176E - 02	0.2834E - 02	0.3233E - 02
85	-0.1783E + 03	0.3669E - 05	0.5295E - 03	0.6583E - 03	0.6583E - 03
84	-0.1783E + 03	-0.1613E - 07	0.1336E - 04	0.1831E - 04	0.1831E - 04
81	-0.2134E + 03	-0.3332E - 03	0.4104E - 03	0.9463E - 03	0.1334E - 02
80	-0.2564E + 03	-0.2121E - 03	0.1986E - 03	0.5018E - 03	0.7040E - 03
79	-0.2994E + 03	-0.1180E - 03	0.1051E - 03	0.2651E - 03	0.4512E - 03
78	-0.3414E + 03	-0.8606E - 04	0.5279E - 04	0.2056E - 03	0.3998E - 03
77	-0.3765E + 03	-0.4833E - 04	0.1871E - 04	0.8964E - 04	0.1225E - 03

class y traffic is $f_y = 0.102$. In this case there are a total of 100 states of which 24 are overload states, ($\omega = 76$), giving the same number of stable modes. Their eigenvalues are shown in the second column of Table 1, using the numbering convention of Section 2. (To exhibit the dominant eigenvalues, the entries are listed in decreasing order of z_n . Note that this does not correspond to decreasing order of n .) The remaining columns compare the exact values of the coefficients b_n in the expression for $G(x)$ (Equation (38)) with the bounds \underline{b}_n , \bar{b}_n , \hat{b}_n defined in (46), (45), and (47) respectively.

Note that there are only three coefficients, b_{83} , b_{92} , b_{100} , of significant magnitude, and the bounds on each are fairly tight. While the bounds on the smaller coefficients are poor in a relative sense, they are still quite good in an absolute sense. One can use the coefficient bounds \underline{b}_n and \bar{b}_n of Table 1 to obtain (approximate) bounds on $G(x)$. Figure 7 shows comparisons of $\bar{G}_{10}(x)$ and $G_{10}(x)$ as defined in (48) and (49) respectively, with the exact value $G(x)$ for the system of Figure 5. Note that these expressions bound G_{10} in (40) using upper and lower bounds on the coefficients of each exponential. Thus Figure 7 illustrates the combined effects of using a subset of the modes and of bounding their coefficients. The expressions do not require the calculation of any system modes beyond those contained in the summation. Furthermore, they do not require solution of any linear equations for the coefficients.

To illustrate the computational techniques in the setting of a larger system we treat the case $N^x = 29$, $N^y = 29$ (900 states) in Figure 8. This corresponds to a tripling of

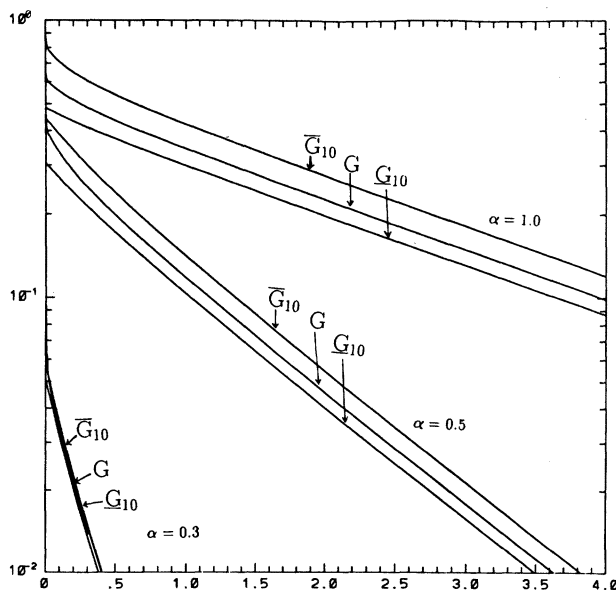


Figure 7. Probability of overflow G and its approximations \bar{G}_{10} and G_{10} vs. buffer size for $N_x = N_y = 9$. For $\alpha = 1, 0.5$, and 0.3 , the traffic intensity ρ is set to $0.91, 0.87$, and 0.74 , respectively

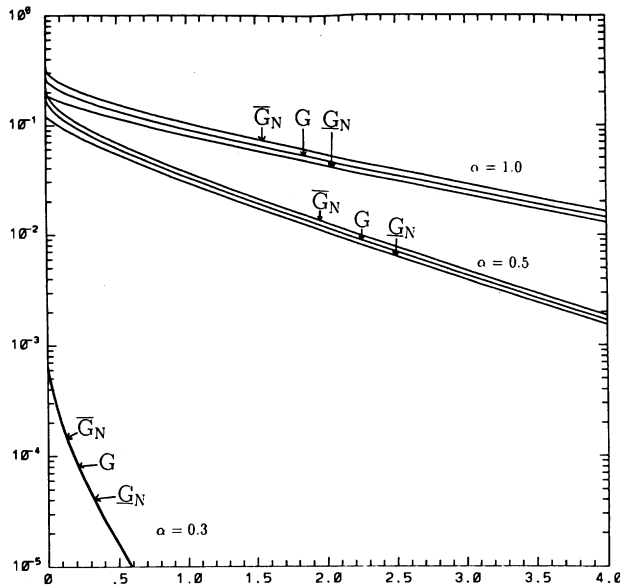


Figure 8. Probability of overflow G and its approximations \bar{G}_N and \tilde{G}_N vs buffer size for $N_x = N_y = 29$. For $\alpha = 1, 0.5$, and 0.3 , the traffic intensity ρ is set to $0.91, 0.87$, and 0.74 , respectively

the size of the original system, with an increase of a factor of 9 in the number of states. The capacity C is adjusted to maintain approximately the same traffic intensities as in the smaller system of Figure 7. The curves show $\bar{G}_N(x)$ and $\tilde{G}_N(x)$ as defined in (48) and (49) for comparison with the exact value of $G(x)$. Note that this corresponds to using all of the modes with exact coefficients replaced by upper and lower bounds respectively. The results give very tight bounds on the exact solution, especially in the case $\alpha = 0.3$. It can be seen that the buffer occupancy for this system decreases markedly as system size increases while keeping the traffic intensity fixed. One would expect this as a consequence of the laws of large numbers.

We observe from these examples that the algorithm is numerically stable, performs efficiently (computation time was relatively short), and the coefficient bounds are well suited to engineering approximations.

6. Conclusions

This work has examined a broad class of buffered information-handling systems, modeled as continuous flow processes modulated by an underlying reversible Markov process. We showed that if the modulating process is separable, considerable reduction in computational complexity is possible using a decomposition method. This avoids problems of 'state explosion' and numerical inaccuracy in very large systems. Further reductions in computational complexity were shown to be possible using simple bounds on coefficients in the expression for the buffer

occupancy distribution. The behavior of the various computational algorithms was illustrated in an example wherein the performance of a multiplexer involving channel sharing was computed. While the largest system illustrated have had 900 states, our computational experiments give us confidence that much larger systems can be handled.

We have been extending this work in a number of directions. As mentioned in Section 1, when the system represented by (5) is composed of the superposition of a large number of identical subsystems the total number of states can be significantly reduced by aggregating the states of the global system. Work is in progress on adapting the results presented here to this special case. Work is also partly completed on carrying over the essential features of these results to the case of point process models: i.e., Markov modulated Poisson processes [5]. In this case, the models are similar to those treated in [18], with the additional features of reversibility and separability. An objective of this effort is to show how certain computational problems associated with matrix-geometric methods can be avoided using decomposition.

Appendix

A.1. *Proof of Equations (18)–(20).* Because the matrices D and \tilde{M} are both Hermitian, it can be shown that (17) has real eigenvalues and a set of N (real) linearly independent eigenvectors. The proof is based on simultaneous diagonalization of two quadratic forms. (See, for example, [7], p. 106 or the appendix of [15].)

To prove (18) consider two pairs $\{z_j, \tilde{\varphi}_j\}$ and $\{z_i, \tilde{\varphi}_i\}$ satisfying (17). Taking inner products on each side of (17) we have

$$z_j(\tilde{\varphi}_i, D\tilde{\varphi}_j) = (\tilde{\varphi}_i, \tilde{M}\tilde{\varphi}_j) = (\tilde{\varphi}_j, \tilde{M}\tilde{\varphi}_i) = z_i(\tilde{\varphi}_j, D\tilde{\varphi}_i) = z_i(\tilde{\varphi}_i, D\tilde{\varphi}_j)$$

or

$$(z_i - z_j)(\tilde{\varphi}_i, D\tilde{\varphi}_j) = 0,$$

which implies (18) when the two eigenvalues are distinct. In the case of a set of identical eigenvalues, any basis for the subspace spanned by their eigenvectors can be orthogonalized with respect to the matrix D using the Gram-Schmidt orthogonalization procedure.

To prove (19) we expand the left side to

$$(\tilde{\varphi}_1, D\tilde{\varphi}_1) = (\tilde{\varphi}_1, R\tilde{\varphi}_1) - C(\tilde{\varphi}_1, \tilde{\varphi}_1).$$

From the definition of E and the fact that $\varphi_1 \sim \pi$ we have

$$(\tilde{\varphi}_1, R\tilde{\varphi}_1) = \bar{r}(\tilde{\varphi}_1, \tilde{\varphi}_1),$$

from which (19) follows directly.

Finally, to prove (20) we take inner products on each side of (17) using the same eigenvectors, to give

$$z_i(\tilde{\varphi}_i, \mathbf{D}\tilde{\varphi}_i) = (\tilde{\varphi}_i, \tilde{\mathbf{M}}\tilde{\varphi}_i)$$

from which (20) follows.

A.2. Properties of the eigenvalues. To prove the properties of the eigenvalues stated in Section 2.1, we need to relate the eigenvalues and eigenvectors of (14) to those of the matrix $\mathbf{A}(z)$ defined in (22). We first prove that for $z \neq 0$, $\det[\mathbf{D}z - \mathbf{M}] = 0$ if and only if $\det[\mathbf{C}\mathbf{I} - \mathbf{A}(z)] = 0$. The equation satisfied by the eigenvalues $\{z_n\}$ can be written

$$\det[\mathbf{D}z - \mathbf{M}] = \det[\mathbf{R}z - \mathbf{C}z\mathbf{I} - \mathbf{M}] = (-z)^N \det[\mathbf{C}\mathbf{I} - \mathbf{A}(z)] = 0.$$

Thus, there must be exactly $(N - 1)$ non-zero values of z for which \mathbf{C} is an eigenvalue of $\mathbf{A}(z)$, corresponding to the $(N - 1)$ non-zero eigenvalues of (14).

Now, to exhibit the correspondence between eigenvectors, let $g_n(z)$ denote the n th eigenvalue of $\mathbf{A}(z)$, with eigenvector $\varphi(n; z)$; i.e.

$$(74) \quad \mathbf{A}(z)\varphi(n; z) = g_n(z)\varphi(n; z) \quad n = 1, 2, \dots, N, \quad z \neq 0,$$

which can be rearranged to

$$(75) \quad [\mathbf{R} - g_n(z)\mathbf{I}]z\varphi(n; z) = \mathbf{M}\varphi(n; z).$$

Suppose there exists a non-zero z_n such that

$$(76) \quad \mathbf{C} = g_n(z_n).$$

Then (75) becomes

$$\mathbf{D}z_n\varphi(n; z_n) = \mathbf{M}\varphi(n; z_n)$$

implying that z_n is a solution of (14) with eigenvector

$$\varphi_n = \varphi(n; z_n).$$

We now show that the functions $g_n(z)$ can be defined so that for $n = 2, \dots, N$, (76) has a unique solution z_n for any $C > \bar{r}$ and C not equal to any of the $r(n)$. The following result will be used. (See [20], p. 44.)

Let $\mathbf{H}(z)$ be an $N \times N$ Hermitian matrix, continuously differentiable with respect to the parameter z , and let its derivative $\mathbf{H}'(z)$ be Hermitian. Then there exist real functions $\mu(n; z)$, $n = 1, 2, \dots, N$, continuously differentiable in z such that

$$\mathbf{H}\xi(n; z) = \mu(n; z)\xi(n; z)$$

where $\{\xi(n; z)\}$ is a properly chosen orthonormal system of vector functions. At every point at which $\mathbf{H}(z)$ has distinct eigenvalues, their derivatives are given by

$$(77) \quad \mu'(n; z) = (\xi(n; z), \mathbf{H}'(z)\xi(n; z)).$$

In the case at hand, we let

$$H(z) = z\tilde{A}(z) = zE^{-1}AE = Rz - \tilde{M}$$

which is Hermitian, so that $\mu(n; z)$ is an eigenvalue of $z\tilde{A}(z)$ for $n = 1, 2, \dots, N$, $\{\xi(n; z)\}$ is an orthonormal set of eigenvectors of $z\tilde{A}(z)$ as well as $\tilde{A}(z)$ and $\mu'(n; z) = (\xi(n; z), R\xi(n; z))$.

To define an ordering of the $\{\mu(n; z)\}$, let

$$0 = \lambda_1 < \lambda_2 \leq \dots \leq \lambda_N$$

be the eigenvalues of $-\tilde{M}$, and let $\mu(n; 0) = \lambda_n$. First, consider the case where $z\tilde{A}(z)$ has distinct eigenvalues for all z . Then by continuity, the ascending order of the $\{\mu(n; z)\}$ must be maintained for all z . Now let

$$(78) \quad g_n(z) = \begin{cases} \mu(n; z)/z, & z > 0, \\ \mu(N - n + 1; z)/z, & z < 0. \end{cases}$$

For each $z \neq 0$, $g_n(z)$ is the n th eigenvalue of $\tilde{A}(z)$ (in ascending order), with eigenvector

$$(79) \quad \tilde{\varphi}(n; z) = E^{-1}\varphi(n; z) = \begin{cases} \xi(n; z), & z > 0, \\ \xi(N - n + 1; z), & z < 0. \end{cases}$$

It follows that each $g_n(z)$ is continuously differentiable for all $z \neq 0$ with

$$(80) \quad g'_n(z) = (\tilde{\varphi}(n; z), \tilde{M}\tilde{\varphi}(n; z))/z^2 \leq 0.$$

Since \tilde{M} is negative semi-definite with null-space spanned by $\tilde{\varphi}_1$, $g'_n(z) = 0$ if and only if $\tilde{\varphi}(n; z) = c\tilde{\varphi}_1$ for some constant c . But this implies

$$\tilde{A}\tilde{\varphi}_1 = (R - \tilde{M}/z)\tilde{\varphi}_1 = R\tilde{\varphi}_1 = g_n(z)\tilde{\varphi}_1.$$

Since all elements of $\tilde{\varphi}_1$ are non-zero, this relation can only hold if R is a scalar matrix, contradicting the inequalities in (8). Thus we conclude that the inequality in (80) is strict for all $z \neq 0$.

It follows directly from the limiting values of $A(z)$ at 0 and ∞ , and the continuity and ordering of the set $\{g_n(z)\}$ that

$$\begin{aligned} \lim_{z \rightarrow 0^+} g_n(z) &= +\infty, & n = 2, 3, \dots, N, \\ \lim_{z \rightarrow 0^-} g_n(z) &= -\infty, & n = 1, 2, \dots, N - 1, \\ \lim_{z \rightarrow \pm\infty} g_n(z) &= r(n), & n = 1, 2, \dots, N. \end{aligned}$$

To find the remaining limits we use l'Hospital's rule to give

$$\lim_{z \rightarrow 0^+} g_1(z) = \lim_{z \rightarrow 0^-} g_N(z) = \mu'(1; 0) = (\tilde{\varphi}(1; 0), R\tilde{\varphi}(1; 0)) = (\tilde{\varphi}_1, R\tilde{\varphi}_1) = \bar{r}.$$

The above development holds with minor variations for the case where $z\bar{A}(z)$ has repeated eigenvalues for some or all values of z . In this case, some of the curves $\mu(n; z)$ may intersect, or coincide over the entire range of z . If they intersect, the ascending order of the $\mu(n; z)$ may not be maintained for all z , which means that with the mapping from μ to g defined in (78), the limits at ∞ of $g_n(z)$ defined in (81) may occur in a different order; e.g.

$$\lim_{z \rightarrow +\infty} g_n(z) = r(i), \quad i \neq n.$$

To maintain the correct order of the limiting values, the mapping in (78) must then be replaced by

$$(81) \quad g_n(z) = \begin{cases} \mu(i(n); z)/z, & z > 0, \\ \mu(j(n); z)/z, & z < 0 \end{cases}$$

where $i(n)$ and $j(n)$ are appropriate permutations of the integers, $1, 2, \dots, N$. Continuous differentiability of the eigenvalues still holds in this case, but one must be more careful about defining the derivatives. (See [20].)

The properties of the family $\{g_n(z)\}$ derived above and illustrated in Figure 4 show that for $n = 2, 3, \dots, N - 1$, $g_n(z)$ is a one-one mapping of $(0, \infty)$ onto $(r(n), \infty)$ and $(-\infty, 0)$ onto $(-\infty, r(n))$, while $g_N(z)$ is a one-one mapping of $(0, \infty)$ onto $(r(N), \infty)$ and $(-\infty, 0)$ onto $(\bar{r}, r(N))$. We therefore conclude that the equation

$$C = g_n(z)$$

has a unique solution, z_n , for any $C > \bar{r}$ and not equal to any of the $r(n)$. Furthermore, with ω defined in (8), it is clear from the form of $g_n(z)$ that

$$\begin{aligned} z_n > 0, & \quad 2 \leq n \leq \omega \\ < 0, & \quad \omega < n \leq N. \end{aligned}$$

A.3. *Proof of Equation (44)*. Substituting (19) into (31) for the case $i = 1$, we have

$$1 = (\tilde{\varphi}_1, DE^{-1}f)/(\bar{r} - C)$$

or

$$\sum_s \tilde{\varphi}_1(s)e^{-1}(s)d(s)f(s) = \sum_s -d(s)f(s) = C(1 - \rho).$$

A.4. *Kronecker product and sum: definitions and properties*. The Kronecker product $A \otimes B$ of the matrix A of dimension $p \times q$ and the matrix B of dimension $m \times n$ is the matrix of dimension $pm \times qn$ obtained by replacing each element a_{ij} of the matrix A by the full matrix $a_{ij}B$. (See for example, [2].)

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

$$A \oplus B = A \otimes I_m + I_n \otimes B$$

where I_m and 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 [3], [2], [8], are used in this paper. All matrices (vectors) are assumed to be of appropriate order.

1. $(A \otimes B)(C \otimes D) = (AC) \otimes (BD)$
2. $(A \otimes B, C \otimes D) = (A, C)(B, D)$
3. 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 $j = 1, 2, \dots, m$.

A.5. Decomposition. In this section we prove relations (57)–(66). It suffices to show that the relations stated in (21) and (22) and proved in Section A.2 are expressible in the decomposed form of (58), (59), and (60). Recall from Section A.2 that for non-zero z_n , $\{z_n, \varphi_n\}$ is a solution of (14) if and only if it satisfies

$$(82) \quad C\varphi_n = A(z_n)\varphi_n.$$

Now the separability conditions imply that

$$(83) \quad A(z) = A^{(1)}(z) \oplus A^{(2)}(z) \oplus \dots \oplus A^{(K)}(z).$$

Let $g(z)$ denote an eigenvalue of $A(z)$, with eigenvector $\varphi(z)$. Then from Property 3 in Section A.4 it follows that

$$(84) \quad g(z) = \sum_k g^{(k)}(i^{(k)}; z)$$

for some set of integers $i^{(1)}, i^{(2)}, \dots, i^{(K)}$, where $g^{(k)}(i^{(k)}; z)$ represents the $i^{(k)}$ th eigenvalue of $A^{(k)}(z)$ in asymptotically ascending order. Furthermore,

$$(85) \quad \varphi(z) = \varphi^{(1)}(i^{(1)}; z) \otimes \varphi^{(2)}(i^{(2)}; z) \otimes \dots \otimes \varphi^{(K)}(i^{(K)}; z)$$

where

$$(86) \quad g^{(k)}(i^{(k)}; z)\varphi^{(k)}(i^{(k)}; z) = A^{(k)}(z)\varphi^{(k)}(i^{(k)}; z).$$

Conversely, for any set $\{i^{(k)}\}$ and any $z \neq 0$ a pair $\{g(z), \varphi(z)\}$ satisfying (84), (85), (86) must satisfy (82) with $C = g(z)$.

If the functions $g(z)$ obtained in this way are indexed lexicographically to define the set $\{g_n(z)\}_{n=2}^N$, then these functions are equivalent to those defined for the global system in (21). In particular, each function $g_n(z)$, $n = 2, 3, \dots, N - 1$ is a one-one mapping of $(0, \infty)$ onto $(r(s_n), \infty)$ and $(-\infty, 0)$ onto $(-\infty, r(s_n))$, while $g_N(z)$ is a one-one mapping of $(0, \infty)$ onto $(r(s_N), \infty)$ and $(-\infty, 0)$ onto $(\bar{r}, r(s_N))$. It therefore follows from Section A.2 that all non-zero eigenvalues of (14) are obtained as solutions of (57) with the corresponding eigenvectors expressed as the Kronecker products of (60). This proves (57)–(60). Properties (62)–(65) follow

directly from Section A.2, since each subsystem satisfies the same reversibility assumption as the global system. The identifications of the signs of the eigenvalues in (66) follows directly from the asymptotic behavior of $g_n(z)$, inherited from its component functions on the right-hand side of (58).

Acknowledgments

The authors acknowledge with thanks, helpful discussions with Christian Huitema and Philippe Nain of INRIA, and computational help from Clark Wang of Columbia. Thanks are also due to D. Mitra for making a pre-publication copy of [15] available.

References

- [1] ANICK, D., MITRA, D. AND SONDDHI, M. M. (1982) Stochastic theory of a data-handling system with multiple sources. *Bell System Tech. J.* **61**, 1871–1894.
- [2] BELLMAN, R. (1960) *Introduction to Matrix Analysis*. McGraw-Hill, New York.
- [3] BREWER, J. W. (1978) Kronecker products and matrix calculus in system theory. *IEEE Trans. Circuits Syst.* **25**, 772–781.
- [4] BURMAN, D. Y. AND SMITH, D. R. (1986) An asymptotic analysis of a queueing system with markov-modulated arrivals. *Operat. Res.* **34**, 05–119.
- [5] ELWALID, A. E., MITRA, D. AND STERN, T. E. (1990) A theory of statistical multiplexing of markovian sources: spectral expansions and algorithms. In *Numerical Solution of Markov Chains*, ed. W. J. Stewart, Marcel Dekker, New York.
- [6] FISCHER, M. J. (1979) Data performance in a system where data packets are transmitted during voice silent periods—single channel case. *IEEE Trans. Comm.* **27**, 1371–1375.
- [7] FRANKLIN, J. N. (1968) *Matrix Theory*. Prentice Hall, Englewood Cliffs, NJ.
- [8] GRAHAM, A. (1981) *Kronecker Products and Matrix Calculus with Applications*. Ellis Horwood, Chichester.
- [9] KEILSON, J. AND RAO, S. S. (1970) A process with chain-dependent growth rate. *J. Appl. Prob.* **7**, 699–711.
- [10] KEILSON, J. AND RAO, S. S. (1971) A process with chain-dependent growth rate. part II: the ruin and ergodic problems. *Adv. Appl. Prob.* **3**, 315–338.
- [11] KEILSON, J. (1979) *Markov Chain Models—Rarity and Exponentiality*. Springer-Verlag, New York.
- [12] KELLY, F. P. (1979) *Reversibility and Stochastic Networks*. Wiley, Chichester.
- [13] KOSTEN, L. (1984) Stochastic theory of data-handling systems with groups of multiple sources. In *Performance of Computer-Communication Systems*, ed. H. Ruding and W. Bux, Elsevier, Amsterdam, 321–331.
- [14] KOSTEN, L. (1986) Liquid models for a type of information buffer problems. *Delft Progress Report* **11**, 71–86.
- [15] MITRA, D. (1988) Stochastic theory of a fluid model of producers and consumers coupled by a buffer. *Adv. Appl. Prob.* **20**, 646–676.
- [16] MITRANI, I. (1968) A many server queue with service interruptions. *Operat. Res.* **16**, 628–638.
- [17] MORRISON, J. (1989) Asymptotic analysis of a data-handling system with many sources, *SIAM J. Appl. Math.* **49**, 617–637.
- [18] NEUTS, M. (1981) *Matrix Geometric Solutions in Stochastic Systems*. Johns Hopkins University Press, Baltimore.
- [19] NEUTS, M. (1986) The caudal characteristic curve of queues. *Adv. Appl. Prob.* **18**, 221–254.

- [20] RELICH, F. (1969) *Perturbation Theory of Eigenvalue Problems*. Gordon and Breach, New York.
- [21] SCHWARTZ, M. (1987) *Telecommunication Networks, Protocols, Modeling and Analysis*. Addison-Wesley, Reading, Mass.
- [22] STERN, T. E. (1983) A queueing analysis of packet voice. *Proc. GLOBECOM '83*, 71–76.
- [23] STERN, T. E. AND ELWALID, A. I. (1989) Analysis of separable markov-modulated rate models for information-handling systems. CTR technical report No. 164–89-43, Columbia University.
- [24] WEISS, A. (1986) A new technique for analyzing large traffic systems. *Adv. Appl. Prob.* **18**, 506–532.