

Robust Dynamic Admission Control for Unified Cell and Call QoS in Statistical Multiplexers

Debasis Mitra, *Member, IEEE*, Martin I. Reiman, and Jie Wang

Abstract—The design of connection admission control (CAC) for a simple Markovian model of a multiservice statistical multiplexer is considered. The paper begins by laying the foundation through several fundamental analytic concepts, such as a semi-Markov decision process formulation of the design problem and time scale decomposition, before progressively leading up to real-world requirements, like robustness and simplicity of design. Several numerical illustrations are given. The salient contributions of the paper are as follows. 1) A unified treatment of multiclass cell and call QoS. 2) A CAC design which is robust, fair, and efficient. 3) Simplicity in the CAC design, together with an evaluation of the tradeoff with performance. 4) An analytic technique for computing the feasibility region in the space of call arrival rates where some control exists to satisfy QoS. 5) The discovery of near linearity of the boundary of the feasible region, which is then used to decompose the design problem. 6) A unified treatment of aggressive and conservative forms of CAC, the latter being conventional and the former yielding better call level performance. 7) An effective bandwidth definition based on the aggressive form of CAC, which influences the CAC design. 8) Demonstration of the beneficial impact on performance of cell level control. 9) An asymptotic theory of the joint behavior of cell loss and call blocking. 10) A rigorous development of time scale decomposition. 11) A numerical evaluation of the accuracy of the notion of nearly completely decomposable Markov chains.

I. INTRODUCTION

THIS paper considers the design of connection admission control (CAC) for a simple Markovian model of a multiservice statistical multiplexer. The dominant characteristic of the CAC considered here is that it addresses both cell and call Quality of Service (QoS) issues, which is different from the conventional focus on the cell level only. The paper begins by dealing rigorously with several fundamental analytic concepts, such as time scale decomposition between the burst and call level and a Semi-Markov Decision Process formulation of the call admission control problem. This is accompanied by several numerical examples. The paper progressively leads up to real world requirements, such as robustness and simplicity of the CAC design. These are meta-problems for which there are no simple problem statements, yet good designs are usually easy to identify. For the analytically intractable design issues, in the latter part of the paper, we rely on a combination of sound fundamental bases, empiricism, and validation.

Manuscript received December 1997; revised February 1998. Parts of this paper were presented at the 34th IEEE Conference on Decision and Control, 1995, and the 15th International Teletraffic Congress, Washington, DC, 1997.

D. Mitra and M. I. Reiman are with Bell Laboratories, Lucent Technologies, Murray Hill, NJ 07947 USA.

J. Wang is with AT&T Labs, Holmdel, NJ 07733 USA.

Publisher Item Identifier S 0733-8716(98)04150-X.

First, why combine cell and call issues in CAC design? Because, ultimately, the service provider will dimension facilities to meet both cell and call level QoS requirements (where the call level QoS requirement is that the blocking probability of each class be below some level), and ignoring call level QoS in CAC design will lead to overengineered systems, as we now argue. Recall that the principal goal of connection admission control in a broadband system is to simultaneously maintain the QoS for several traffic streams with different characteristics. There is a substantial literature in this area (see [13], [21], [11], [23] and references therein). With very few exceptions, such as [4] and [11], the published work has focused on cell level QoS, such as cell loss ratio, cell delay and cell delay variation. Most CAC schemes are based on the well-known concept of effective bandwidth, which has been studied extensively (see [3], [6], [10], [14]). The formulation of effective bandwidth in the literature is, in turn, mainly based on cell level QoS. The basic call admission principle is tacitly *complete sharing* (CS), i.e., a new call is admitted if by doing so the previously agreed cell level QoS of this and other calls already in progress will not be violated. However, in order to take advantage of statistical multiplexing, it is necessary to carry traffic with different statistical characteristics on the same network and CS can cause a significant difference in call blocking probabilities between sources with different bandwidth requirements. Thus, in order to meet the call level QoS requirement for higher bandwidth calls under CS, the facility size will be such that low bandwidth calls will see call blocking probabilities substantially lower than required. Hence, in order to obtain efficient resource sharing, it is necessary to consider not only the cell level QoS but also the call level QoS when designing the CAC. This leads naturally to call admission control schemes that are not CS. In this paper we propose two ways to take call level QoS into consideration when studying call admission control. First, we introduce a new measure of bandwidth requirement for connections that takes into account both cell and call level QoS requirements. Second, based on the new bandwidth requirement, which we call (yet another) *effective bandwidth*, we propose a new method for call admission control, which guarantees not only cell level QoS, but also regulates traffic at the call level to meet the call level QoS requirements in a robust manner.

There is a particular notion of robustness employed here. For each service or traffic class we have a corresponding engineered load, with the complete set of such loads forming part of the input to the CAC design problem. In reality, of course, the offered load for a class may be more or less than

its engineered load, and we use the terms “overload” and “underload,” respectively, to distinguish the loading status of the class. Our notion of robustness and fairness requires the underloaded classes to be protected, i.e., their received QoS is only marginally affected when the total load exceeds the engineered level. This protective or isolating feature pulls the design in the direction of *complete partitioning*, which is at the other end of the spectrum from complete sharing. The CAC design we propose here is based on the resource-sharing technique of *virtual partitioning* (VP) [18], [5]. This technique aims to achieve controlled sharing, which strikes a balance between unrestricted sharing and isolation. Instead of each class of traffic having a fixed priority, as in traditional trunk reservation ([1], [15]), in virtual partitioning the priorities depend on the state of the system.

In the first part of the paper we examine several topics that are important for providing the bases of our CAC design. We begin in Section II by introducing our model, which consists of a single bufferless link with multiple call classes. The call classes may have quite different peak rates and burstiness properties, and each call behaves as an on-off fluid source while in the system. We then consider the following optimization problem. Given some maximum tolerable cell loss, devise a call admission control procedure that provides acceptable cell loss and maximizes the revenue due to carried traffic. We consider the optimization problem both with and without constraints on the call blocking probabilities. The optimization problem is formulated as a semi-Markov decision process (SMDP) with constraints.

For any system of realistic size, the optimization problem of Section II is so computationally intensive that, in effect, it is numerically impractical. This leads us to seek some method that would enable us to reduce the computational burden. There is a natural time scale decomposition in our system that arises due to the large disparity between the lifetime of a burst and the lifetime of a call. This time scale decomposition allows us to apply the notion of nearly completely decomposable (NCD) Markov chains ([2], [11], [9]) to our system. The NCD approximation reduces the dimension of the state, and the reduced state optimization is again formulated as an SMDP in Section III. The reduced state optimization problem is numerically feasible for systems of realistic size. Two related versions of the problem are considered: conservative and aggressive. In the former approach we require that the cell loss constraints are satisfied for every state, i.e., the number of calls of each class in progress. In the aggressive approach, we only require that the cell loss constraints be satisfied on a long run average basis. Next, also in Section III, we describe how cell level control, i.e., selective cell discarding based on the cell class, can be used in a two-class problem to reduce two loss constraints to one constraint. Accompanying numerical results show that the gain in capacity from cell level control is significant, especially if the required cell loss ratios are quite different. (Although the cell level control does not play an explicit role in the design of our simple CAC, in formulating the constraints on cell level behavior we implicitly assume that an optimal cell level control is being used.) The final topic considered in Section III is the accuracy

of the NCD approximation. Specifically, we give results and insights from a numerical investigation on the solution of the SMDP problem, with and without the NCD approximation. The results indicate that the NCD approximation works well when the average number of on-off cycles during a call's holding time is about 100 or larger.

Section IV is a short, but important, bridge between the first half of the paper, which is dominated by the SMDP and NCD concepts, and the second half, where the focus is on simple and robust CAC designs. In this section we ask if, given the statistical parameters of the burst and call level behaviors of each of the two classes, and the cell and call QoS parameters, there exists a CAC such that the QoS constraints are satisfied? This is a feasibility question and, not surprisingly, it is determined by the feasibility of a linear program (LP). The answer to the question is exhibited as a region, called the “feasibility region,” in λ -space, i.e., where the coordinates are the arrival rates, λ_i , of calls. An important discovery is that the boundary of the feasible region exhibits near-linearity. Exploitation of this feature allows multiple class problems to be treated in the next section as multiple single class problems. Another feature of the numerical results of this section is that the feasible region for the aggressive approach is significantly larger than for the conservative approach. With this as motivation, we consider only the aggressive approach in the remainder of the paper.

Starting with Section V we consider the realization of simple and robust CAC for feasible λ . Section V considers a single class system and proposes the new effective bandwidth. Importantly, the effective bandwidth as calculated here is less conservative than traditionally calculated values which ignore call dynamics. This is because in the latter the worst case call configuration dominates, while here the call distribution is taken into account. An asymptotic analysis of the system is given in Section VI for the scaling in which the call arrival rates and the link bandwidth are simultaneously made large. Although these asymptotics are not used for our CAC design, they provide a valuable insight into the parameter scalings that arise as a consequence of the qualitatively different scalings in the cell and call level QoS constraints. In Section VII we discuss the standard policies of complete sharing and trunk reservation. Virtual partitioning is introduced in Section VIII. In Section IX we present our numerical results on the performances of various call admission control policies, with emphasis on robustness. Our results clearly indicate that our proposed CAC is able to protect the underloaded class in the presence of traffic that deviates from the engineered load, and that neither complete sharing nor trunk reservation provides this protection. New directions are discussed in Section X.

We end this section by summarizing the salient contributions of this paper. 1) A unified treatment of multiclass cell and call QoS in CAC design. 2) The concept of robust CAC design based on a balance of fairness and efficiency. 3) Simplicity of design. The tradeoff between simplicity and performance is highlighted in the numerical work, where it is shown that a small portion of the feasibility region is traded for simplicity. 4) An analytic technique for computing the feasibility region. The technique handles small cell loss probabilities, such as

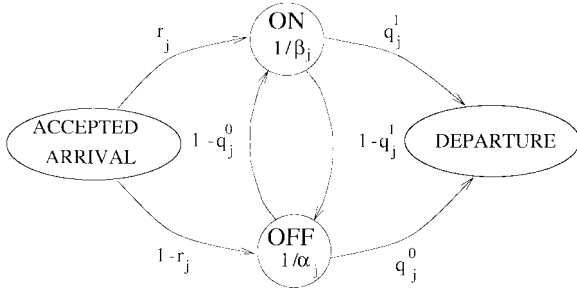


Fig. 1. A Markov chain describing call and burst dynamic.

10^{-9} , which simulations cannot. 5) The discovery of near-linearity of the boundary of the feasibility region. This is an example of the empirical content of the paper. 6) A unified treatment of conservative and aggressive forms of CAC. The conservative form has been conventionally accepted, yet the paper provides solid grounds based on performance for favoring the latter. 7) An effective bandwidth definition based on the aggressive form of CAC. 8) The concept of cell level control and a quantification of the gain in performance. 9) The asymptotic theory of the joint behavior of cell loss and call blocking, where the scaling is for exponentially small cell loss and critical loading at the call level. 10) A rigorous development of time scale decomposition. 11) A numerical evaluation of the accuracy of the NCD technique.

II. CALL ADMISSION CONTROL: A SEMI-MARKOV DECISION PROCESS FORMULATION

A. The Model

We consider the call admission control problem for a statistical multiplexer system that consists of a single link with transmission rate R and no buffer. There are J classes of calls arriving according to independent Poisson processes, and class j calls arrive at rate λ_j , $j = 1, \dots, J$. After being admitted into the system, a class j ($1 \leq j \leq J$) call behaves according to a two-state Markov process depicted in Fig. 1. The two states represent “on” (bursting) and “off.” The initial state is “on” with probability r_j , and “off” with probability $1 - r_j$. The mean holding times in the “on” and “off” states are β_j^{-1} and α_j^{-1} , respectively. A call leaving the “on” (“off”) state has a probability of q_j^1 (q_j^0) of departing from the system. When a class j call is in the “on” state, it generates cells as fluid at rate ν_j . When a class j call is admitted, the system collects a fixed amount of revenue w_j . Note that, in contrast with [8], where a decision theoretic framework is used to deal with unknown parameters, we are assuming that the network has information on the statistical characteristics of the traffic at the burst and call levels.

Based on the above model, we can calculate the average call holding time for class j , $1/\mu_j$, as a time to exit a two-state transient semi-Markov process. The states are 0 (off), and 1 (on), and the transition matrix of the embedded Markov chain is

$$P^{(j)} = \begin{pmatrix} 0 & 1 - q_j^0 \\ 1 - q_j^1 & 0 \end{pmatrix}.$$

The mean residence time in state 0 is $1/\alpha_j$, and the mean residence time in state 1 is $1/\beta_j$. Let $N^{(j)} = (I - P^{(j)})^{-1}$. Then

$$\frac{1}{\mu_j} = (1 - r_j) \left[N_{00}^{(j)} / \alpha_j + N_{01}^{(j)} / \beta_j \right] + r_j \left[N_{10}^{(j)} / \alpha_j + N_{11}^{(j)} / \beta_j \right].$$

A straightforward calculation shows that

$$N^{(j)} = \frac{1}{1 - (1 - q_j^0)(1 - q_j^1)} \begin{pmatrix} 1 & 1 - q_j^0 \\ 1 - q_j^1 & 1 \end{pmatrix}$$

so that

$$\frac{1}{\mu_j} = \frac{\beta_j(1 - r_j q_j^1) + \alpha_j(1 - (1 - r_j)q_j^0)}{\alpha_j \beta_j [1 - (1 - q_j^0)(1 - q_j^1)]}. \quad (1)$$

The fraction of time a class j call spends in the “on” state is

$$p_j = \frac{\alpha_j(1 - (1 - r_j)q_j^0)}{\alpha_j \beta_j [1 - (1 - q_j^0)(1 - q_j^1)]}. \quad (2)$$

When cells are generated at a rate exceeding the link transmission capacity, those cells that cannot be transmitted are lost.

Two types of QoS are of interest: the cell level QoS and the call level QoS. At the cell level, we want the cell loss ratio to be smaller than p_j^{cell} . A commonly used bound is $p_j^{\text{cell}} = 10^{-9}$. The call level QoS is reflected in call blocking probability bounded by p_j^{call} .

B. The Semi-Markov Decision Process (SMDP) and Linear Programming

Let k_j be the number of class j calls in progress, and n_j the number of class j calls in the “on” state. Then for a stationary admission control policy, the process is Markovian and $(\mathbf{k}, \mathbf{n}) = (k_1, \dots, k_J; n_1, \dots, n_J)$ is a state descriptor. When there are n_j class j calls in the “on” state, $j = 1, \dots, J$, the total cell loss rate is $[\sum_{j=1}^J n_j \nu_j - R]^+$, while the cell arrival rate is $\sum_{j=1}^J \lambda_j$.

The first call admission control we consider here is the solution to the following optimization problem: maximize the long run average revenue while satisfying the cell level QoS requirement. We formulate this problem as an SMDP and use Linear Programming to solve it. We also consider the same problem with the addition of call level QoS requirements.

1) *State Space and Action Sets*: Let I denote the state space, which is the union of two sets—the set of call arrival states, and the set of states corresponding to the rest of the events. A call arrival state has the format $(\mathbf{k}, \mathbf{n}, j)$ which corresponds to the arrival of a new class j call that finds \mathbf{k} calls in progress of which \mathbf{n} calls are in the “on” state. The rest of the events correspond to either a call departure, a call being turned “on,” or a call being turned “off.” We use (\mathbf{k}, \mathbf{n}) to denote a state corresponding to those events, which indicates that there are \mathbf{k} calls in progress of which \mathbf{n} calls are in the “on” state *after* the event. In order to obtain a solution

to the SMDP from linear programming, the state space must be finite. Let M be a fixed large positive integer. We thus consider as our state space $I(M)$, where

$$I(M) = \left\{ (\mathbf{k}, \mathbf{n}, j) : k_i, n_i \in \mathbb{Z}^+, n_i \leq k_i, \right. \\ \left. 1 \leq i \leq J; 1 \leq j \leq J; \text{ and } \sum_{i=1}^J k_i \leq M \right\} \\ \cup \left\{ (\mathbf{k}, \mathbf{n}) : k_i, n_i \in \mathbb{Z}^+, n_i \leq k_i, \right. \\ \left. 1 \leq i \leq J; \text{ and } \sum_{i=1}^J k_i \leq M \right\}.$$

For each state in the first set, $(\mathbf{k}, \mathbf{n}, j)$, the corresponding action set is $A(\mathbf{k}, \mathbf{n}, j) = \{0, 1\}$, where 0 denotes rejection and 1 denotes acceptance. For states of the form (\mathbf{k}, \mathbf{n}) , there is no action to be taken, and we set $A(\mathbf{k}, \mathbf{n}) = \{0\}$.

2) *Mean Sojourn Times:* Let $\tau[(\mathbf{k}, \mathbf{n}, j); a]$ denote the average time spent in state $(\mathbf{k}, \mathbf{n}, j)$ until the next decision epoch if action a is taken; and $\tau[(\mathbf{k}, \mathbf{n}); 0]$ denote the average time spent in state (\mathbf{k}, \mathbf{n}) until the next decision epoch. In addition, let $S(\mathbf{k}, \mathbf{n}) = \sum_{j=1}^J \lambda_j + \sum_{j=1}^J n_j \beta_j + \sum_{j=1}^J (k_j - n_j) \alpha_j$. Then

$$\tau((\mathbf{k}, \mathbf{n}, j); a) = \begin{cases} \frac{r_j}{S(\mathbf{k} + \mathbf{e}_j, \mathbf{n} + \mathbf{e}_j)} + \frac{1-r_j}{S(\mathbf{k} + \mathbf{e}_j, \mathbf{n})} & a = 1, \\ \frac{1}{S(\mathbf{k}, \mathbf{n})} & a = 0. \end{cases} \\ \tau[(\mathbf{k}, \mathbf{n}); 0] = \frac{1}{S(\mathbf{k}, \mathbf{n})}.$$

3) *Transition Probabilities:* Let $p[(\mathbf{k}, \mathbf{n}, j), \mathbf{i}, a]$ denote the transition probability from state $(\mathbf{k}, \mathbf{n}, j)$ to \mathbf{i} if action a is chosen, where $a = 0, 1$; and $p[(\mathbf{k}, \mathbf{n}), \mathbf{i}]$ denote the transition probability from state (\mathbf{k}, \mathbf{n}) to \mathbf{i} . Recall that, for state $(\mathbf{k}, \mathbf{n}, j)$, there are \mathbf{k} calls in progress of which \mathbf{n} calls are in the “on” state *before* the arrival event; and for state (\mathbf{k}, \mathbf{n}) there are \mathbf{k} calls in progress of which \mathbf{n} calls are in the “on” state *after* the event.

Recall that r_j is the probability that an accepted class j call starts in “on” state first. Hence, we have for any $k_j \geq n_j \geq 0$ the equations found at the bottom of the page.

4) *Rewards:* The reward can only be earned in state (k, n, a) when action $a = 1$ is taken

$$r[(\mathbf{k}, \mathbf{n}, j); a] = \begin{cases} 0 & a = 0 \\ w_j & a = 1. \end{cases} \\ r[(\mathbf{k}, \mathbf{n}); 0] = 0.$$

5) *Cell Loss and Arrival Rates:* Recall the previously stated cell loss and arrival rates when n_j class j calls are in the “on” state. Let $l_m[(\mathbf{k}, \mathbf{n}, j); a]$ and $g_m[(\mathbf{k}, \mathbf{n}, j); a]$ denote the expected amount of cells of class m lost and arriving until the next decision epoch when in state $(\mathbf{k}, \mathbf{n}, j)$ and action a is taken; $l_m[(\mathbf{k}, \mathbf{n})]$ and $g_m[(\mathbf{k}, \mathbf{n})]$ denote the expected amount of cells of class m calls lost and arriving in state (\mathbf{k}, \mathbf{n}) until the next decision epoch. Then, with no priority

$$p[(\mathbf{k}, \mathbf{n}, j), \mathbf{i}, 1] = \begin{cases} \frac{r_j \lambda_j}{S(\mathbf{k} + \mathbf{e}_j, \mathbf{n} + \mathbf{e}_j)} & \mathbf{i} = (\mathbf{k} + \mathbf{e}_j, \mathbf{n} + \mathbf{e}_j, l) \quad 1 \leq l \leq J \\ \frac{(1-r_j) \lambda_l}{S(\mathbf{k} + \mathbf{e}_j, \mathbf{n})} & \mathbf{i} = (\mathbf{k} + \mathbf{e}_j, \mathbf{n}, l) \quad 1 \leq l \leq J \\ \frac{r_j (1-q_l^1) (n_l + 1_{\{j=l\}}) \beta_l}{S(\mathbf{k} + \mathbf{e}_j, \mathbf{n} + \mathbf{e}_j)} & \mathbf{i} = (\mathbf{k} + \mathbf{e}_j, \mathbf{n} + \mathbf{e}_j - \mathbf{e}_l) \quad 1 \leq l \leq J \\ \frac{(1-r_j) (1-q_l^1) n_l \beta_l}{S(\mathbf{k} + \mathbf{e}_j, \mathbf{n})} & \mathbf{i} = (\mathbf{k} + \mathbf{e}_j, \mathbf{n} - \mathbf{e}_l) \quad 1 \leq l \leq J \\ \frac{r_j q_l^1 (n_l + 1_{\{j=l\}}) \beta_l}{S(\mathbf{k} + \mathbf{e}_j, \mathbf{n} + \mathbf{e}_j)} & \mathbf{i} = (\mathbf{k} + \mathbf{e}_j - \mathbf{e}_l, \mathbf{n} + \mathbf{e}_j - \mathbf{e}_l) \quad 1 \leq l \leq J \\ \frac{(1-r_j) q_l^1 n_l \beta_l}{S(\mathbf{k} + \mathbf{e}_j, \mathbf{n})} & \mathbf{i} = (\mathbf{k} + \mathbf{e}_j - \mathbf{e}_l, \mathbf{n} - \mathbf{e}_l) \quad 1 \leq l \leq J \\ \frac{r_j (1-q_l^0) (k_l - n_l) \alpha_l}{S(\mathbf{k} + \mathbf{e}_j, \mathbf{n} + \mathbf{e}_j)} & \mathbf{i} = (\mathbf{k} + \mathbf{e}_j, \mathbf{n} + \mathbf{e}_j + \mathbf{e}_l) \\ \frac{(1-r_j) (1-q_l^0) (k_l - n_l + 1_{\{j=1\}}) \alpha_l}{S(\mathbf{k} + \mathbf{e}_j, \mathbf{n})} & \mathbf{i} = (\mathbf{k} + \mathbf{e}_j, \mathbf{n} + \mathbf{e}_l) \\ \frac{r_j q_l^0 (k_l - n_l) \alpha_l}{S(\mathbf{k} + \mathbf{e}_j, \mathbf{n} + \mathbf{e}_j)} & \mathbf{i} = (\mathbf{k} + \mathbf{e}_j - \mathbf{e}_l, \mathbf{n} + \mathbf{e}_j) \\ \frac{(1-r_j) q_l^0 (k_l - n_l + 1_{\{j=1\}}) \alpha_l}{S(\mathbf{k} + \mathbf{e}_j, \mathbf{n})} & \mathbf{i} = (\mathbf{k} + \mathbf{e}_j - \mathbf{e}_l, \mathbf{n}) \\ 0 & \text{otherwise.} \end{cases} \\ p[(\mathbf{k}, \mathbf{n}, j), \mathbf{i}, 0] = p[(\mathbf{k}, \mathbf{n}), \mathbf{i}] = \begin{cases} \frac{\lambda_j}{S(\mathbf{k}, \mathbf{n})} & \mathbf{i} = (\mathbf{k}, \mathbf{n}, j) \\ \frac{q_j^1 n_j \beta_j}{S(\mathbf{k}, \mathbf{n})} & \mathbf{i} = (\mathbf{k}, \mathbf{n} - \mathbf{e}_j) \\ \frac{(1-q_j^1) n_j \beta_j}{S(\mathbf{k}, \mathbf{n})} & \mathbf{i} = (\mathbf{k} - \mathbf{e}_j, \mathbf{n} - \mathbf{e}_j) \\ \frac{q_j^0 (k_j - n_j) \alpha_j}{S(\mathbf{k}, \mathbf{n})} & \mathbf{i} = (\mathbf{k}, \mathbf{n} + \mathbf{e}_j) \\ \frac{(1-q_j^0) (k_j - n_j) \alpha_j}{S(\mathbf{k}, \mathbf{n})} & \mathbf{i} = (\mathbf{k} - \mathbf{e}_j, \mathbf{n}) \\ 0 & \text{otherwise.} \end{cases}$$

cell discarding

$$\begin{aligned}
& \frac{l_m[(\mathbf{k}, \mathbf{n}, j); 1]}{\tau[(\mathbf{k}, \mathbf{n}, j); 1]} \\
&= r_j \left[\sum_{i=1}^J n_i \nu_i + \nu_j - R \right]^+ \frac{n_m \nu_m + 1_{\{m=j\}} \nu_m}{\sum_{i=1}^J n_i \nu_i + \nu_j} \\
&+ (1 - r_j) \left[\sum_{i=1}^J n_i \nu_i - R \right]^+ \frac{n_m \nu_m}{\sum_{i=1}^J n_i \nu_i} \\
& \frac{l_m[(\mathbf{k}, \mathbf{n}, j); 0]}{\tau[(\mathbf{k}, \mathbf{n}, j); 0]} = \frac{l_m[(\mathbf{k}, \mathbf{n}); 0]}{\tau[(\mathbf{k}, \mathbf{n}); 0]} \\
&= \left[\sum_{i=1}^J n_i \nu_i - R \right]^+ \frac{n_m \nu_m}{\sum_{i=1}^J n_i \nu_i} \\
& \frac{g_m[(\mathbf{k}, \mathbf{n}, j); 1]}{\tau[(\mathbf{k}, \mathbf{n}, j); 1]} = n_m \nu_m + 1_{\{m=j\}} r_j \nu_j \\
& \frac{g_m[(\mathbf{k}, \mathbf{n}, j); 0]}{\tau[(\mathbf{k}, \mathbf{n}, j); 0]} = \frac{g_m[(\mathbf{k}, \mathbf{n}); 0]}{\tau[(\mathbf{k}, \mathbf{n}); 0]} = n_m \nu_m.
\end{aligned}$$

6) *Linear Programming Formulation:* Linear programming is a classical technique for solving both Markov decision processes and semi-Markov decision processes, with and without constraints (cf. [22]). With multiple constraints in an SMDP setting, Feinberg [7] shows that it suffices to consider randomized stationary policies for this problem, and the optimal solution subject to class l calls' long-run average cell loss ratio $\leq p_l^{\text{cell}}$ can be obtained from a linear program. The variables in the LP are $z_{\mathbf{i}a}$, $\mathbf{i} \in I(M)$, $a \in A(\mathbf{i})$, where $z_{\mathbf{i}a} \tau(\mathbf{i}, a)$ corresponds to the fraction of time spent in state \mathbf{i} with action a chosen. The LP is

$$\text{maximize} \quad \sum_{\mathbf{i} \in I(M)} \sum_{a \in A(\mathbf{i})} r(\mathbf{i}, a) z_{\mathbf{i}a}$$

subject to

$$\begin{aligned}
& \sum_{a \in A(\mathbf{j})} z_{\mathbf{j}a} - \sum_{\mathbf{i} \in I(M)} \sum_{a \in A(\mathbf{i})} p(\mathbf{i}, \mathbf{j}, a) z_{\mathbf{i}a} = 0, \quad \mathbf{j} \in I(M) \\
& \sum_{\mathbf{i} \in I(M)} \sum_{a \in A(\mathbf{i})} [\ell_j(\mathbf{i}, a) - p_j^{\text{cell}} g_j(\mathbf{i}, a)] z_{\mathbf{i}a} \leq 0, \quad 1 \leq j \leq J \\
& \sum_{\mathbf{i} \in I(M)} \sum_{a \in A(\mathbf{i})} \tau(\mathbf{i}, a) z_{\mathbf{i}a} = 1, \\
& z_{\mathbf{i}a} \geq 0, \quad \mathbf{i} \in I(M), a \in A(\mathbf{i}).
\end{aligned}$$

The above LP is feasible, because $z_{\mathbf{0},0} = [\tau(\mathbf{0}, 0)]^{-1} = \sum_j \lambda_j$, and $z_{\mathbf{i}a} = 0$ otherwise is a feasible solution (corresponding to never allowing any calls to enter). Let $\phi(\mathbf{k}, \mathbf{n}, j)$ denote the probability that action 1 (accept) is chosen in state $(\mathbf{k}, \mathbf{n}, j)$. Given an optimal solution, z , of the LP, we obtain an optimal randomized stationary policy as

$$\phi(\mathbf{k}, \mathbf{n}, j) = \frac{z(\mathbf{k}, \mathbf{n}, j); 1}{z(\mathbf{k}, \mathbf{n}, j); 0 + z(\mathbf{k}, \mathbf{n}, j); 1}$$

if $z(\mathbf{k}, \mathbf{n}, j); 0 + z(\mathbf{k}, \mathbf{n}, j); 1 > 0$, and $\phi(\mathbf{k}, \mathbf{n}, j) = 0$ otherwise.

7) *Both Call and Cell Level QoS Requirements:* We now consider the above optimization problem with the J additional constraints that the call blocking probability for class j should not exceed p_j^{call} , $1 \leq j \leq J$. The treatment in [7] easily enables us to add these constraints to the linear program. It is now possible, however, that the LP will be infeasible: there may be no admission control rule that can simultaneously satisfy all cell level and call level quality of service requirements. (With $J = 2$ and all parameters other than λ_1 and λ_2 fixed, in Section IV we provide a related LP for the reduced state problem introduced in Section III, and numerically determine the set of (λ_1, λ_2) for which the LP is feasible.)

We define, for fixed M , the set of arrival states corresponding to class j calls as

$$\begin{aligned}
I_j^A(M) = & \left\{ (\mathbf{k}, \mathbf{n}, j) : k_i, n_i \in \mathbb{Z}^+, n_i \leq k_i, \right. \\
& \left. 1 \leq i \leq J; \sum_{i=1}^J k_i \leq M \right\}, \quad 1 \leq j \leq J.
\end{aligned}$$

Recall that, in the LP formulation, $\tau(\mathbf{i}, a) z_{\mathbf{i}a}$ corresponds to the fraction of time spent in state \mathbf{i} with action a chosen. Thus, the constraint on class j blocking probability is

$$\sum_{\mathbf{i} \in I_j^A(M)} z_{\mathbf{i}0} - p_j^{\text{call}} \lambda_j \leq 0.$$

The full linear program with both cell and call QoS constraints is thus

$$\text{maximize} \quad \sum_{\mathbf{i} \in I(M)} \sum_{a \in A(\mathbf{i})} r(\mathbf{i}, a) z_{\mathbf{i}a}$$

subject to

$$\begin{aligned}
& \sum_{a \in A(\mathbf{j})} z_{\mathbf{j}a} - \sum_{\mathbf{i} \in I(M)} \sum_{a \in A(\mathbf{i})} p(\mathbf{i}, \mathbf{j}, a) z_{\mathbf{i}a} = 0, \quad \mathbf{j} \in I(M) \\
& \sum_{\mathbf{i} \in I(M)} \sum_{a \in A(\mathbf{i})} [\ell_j(\mathbf{i}, a) - p_j^{\text{cell}} g_j(\mathbf{i}, a)] z_{\mathbf{i}a} \leq 0, \quad 1 \leq j \leq J \\
& \sum_{\mathbf{i} \in I_j^A(M)} z_{\mathbf{i}0} - p_j^{\text{call}} \lambda_j \leq 0, \quad 1 \leq j \leq J \\
& \sum_{\mathbf{i} \in I(M)} \sum_{a \in A(\mathbf{i})} \tau(\mathbf{i}, a) z_{\mathbf{i}a} = 1 \\
& z_{\mathbf{i}a} \geq 0, \quad \mathbf{i} \in I(M), a \in A(\mathbf{i}).
\end{aligned}$$

III. NEARLY COMPLETE DECOMPOSABILITY: IMPLICATIONS AND ACCURACY

The SMDP introduced in Section II is $2J$ dimensional for a J class system. For $J > 1$, this leads to computational problems of excessive size. The notion of NCD Markov chains can be used to reduce this to a J -dimensional problem. (Although the discussion presented here is heuristic, the problem reduction is rigorously justified by Theorem 11 of [2].) Intuitively, this limiting regime is a good approximation when the process describing the number of calls in the “on” state reaches equilibrium between any change in the number

of calls in progress due to call arrival/departure. Hence, when making a call admission decision, the number of calls of each class in progress is important, but the number of calls of each class in the on state is not, because these quantities oscillate too rapidly.

We handle the above idea mathematically as follows. Note that the mean time between call events (call arrival and departure) of class i is $1/(k_i\mu_i + \lambda_i)$. For this quantity to be large relative to the mean burst cycle time, as is required for the NCD limiting regime to hold, it is required that

$$\left(\frac{\alpha_i\beta_i}{\alpha_i + \beta_i} - \lambda_i \right) / \mu_i \gg k_i \quad 1 \leq i \leq J. \quad (3)$$

Now consider a family of system indexed by $\epsilon \rightarrow 0$, in which $\alpha_i(\epsilon) = O(1)$, $\beta_i(\epsilon) = O(1)$, and $r_i(\epsilon) = O(1)$. Assume that

$$q_i^1(\epsilon) = \epsilon q_i^1 \quad \text{and} \quad q_i^0(\epsilon) = \epsilon q_i^0.$$

Then since the average call duration is (see Section II-B)

$$\begin{aligned} \frac{1}{\mu_i(\epsilon)} &= \frac{1 - r_i q_i^1(\epsilon)}{\alpha_i(\epsilon) [1 - (1 - q_i^1(\epsilon))(1 - q_i^0(\epsilon))] \\ &\quad + \frac{1 - (1 - r_i) q_i^0(\epsilon)}{\beta_i(\epsilon) [1 - (1 - q_i^1(\epsilon))(1 - q_i^0(\epsilon))]} \end{aligned}$$

it follows that $\mu_i(\epsilon) = O(\epsilon)$ and

$$\frac{\alpha_i(\epsilon)\beta_i(\epsilon)}{\alpha_i(\epsilon) + \beta_i(\epsilon)} \frac{1}{\mu_i(\epsilon)} = O\left(\frac{1}{\epsilon}\right). \quad (4)$$

Now let us further assume that

$$\frac{\lambda_i(\epsilon)}{\mu_i(\epsilon)} = o\left(\frac{1}{\epsilon}\right). \quad (5)$$

So that, from (4) and (5),

$$\left(\frac{\alpha_i\beta_i}{\alpha_i + \beta_i} - \lambda_i \right) / \mu_i = O\left(\frac{1}{\epsilon}\right). \quad (6)$$

On comparing (3) and (6), we see that the NCD approximation may apply even with a large number of calls in progress. Note that QoS considerations typically prevent the number of calls in progress from being excessively large relative to capacity, as we shall see in our applications. However, whenever the NCD approximation is used, it is necessary to check that the condition in (3) is satisfied. This is the case in the examples treated in the paper.

In the NCD limit, the holding time distribution of a call will be exponential for any on and off time distribution, as long as the memoryless process contained in our model for terminating a call is used. Although it would be possible to incorporate nonexponential call holding times, this would not be entirely straightforward. A SMDP formulation would require that we keep track of the elapsed holding time for each call in progress. (The current phase of each call would be sufficient for a phase-type holding time distribution.) The optimal policy would typically depend in a nonmonotone manner on this elapsed holding time information, making it difficult to implement. An alternative would be to implement a policy that only depends on the number of calls in progress, and ignores the elapsed

holding time information. Such a policy might do well, but it might not.

As $\epsilon \rightarrow 0$, the \mathbf{n} component of the state becomes noise on the time scale where call arrival and departure rates are $O(1)$, and can be ignored for admission control purposes. This part of the state *does* affect the loss rate, so it must be “averaged” properly. For ϵ small, the \mathbf{n} process reaches equilibrium between changes in the \mathbf{k} process. The equilibrium corresponds to fixed \mathbf{k} and is given by the binomial distribution

$$\psi(\mathbf{k}, \mathbf{n}) = \prod_{i=1}^J \binom{k_i}{n_i} p_i^{n_i} (1 - p_i)^{k_i - n_i}.$$

When the total cell arrival rate is $\sum_i n_i \nu_i$, the cell loss rate is $[\sum_i n_i \nu_i - R]^+$. The average class i cell loss rate with \mathbf{k} is thus

$$b_i(\mathbf{k}) = \sum_{n_1=0}^{k_1} \cdots \sum_{n_J=0}^{k_J} \psi(\mathbf{k}, \mathbf{n}) \left[\sum_j n_j \nu_j - R \right]^+ \frac{n_i \nu_i}{\sum_j n_j \nu_j}. \quad (7)$$

This again reflects the fact that there is no priority—the cell loss rate of a class is proportional to its input rate. It seems intuitively clear that the exponential distribution for on and off times plays no essential role in (7). Indeed, the results of [2] apply for on and off times having any phase type distributions.

A. Formulation of the Reduced Problem as an SMDP

1) *States, Actions, and Transition Probabilities:* We now formulate the limit control problem as an SMDP. As in the original problem admission decisions are made upon call arrival. The state space I is again the union of two sets, corresponding to states associated with call arrivals and states associated with call departures. To make this a finite state problem, we may need to place an *a priori* bound on the number of calls that can be accepted in the system. This is taken care of below. Arrival states take the form (\mathbf{k}, j) , with $k_j \geq 0$, and $1 \leq j \leq J$. The state (\mathbf{k}, j) corresponds to an arrival of class j call when there are k_i class i calls in progress, $1 \leq i \leq J$. These are the only states in which decisions need to be made. The set of actions available is $A(\mathbf{k}, j) = \{0, 1\}$, where 0 denotes rejection and 1 denotes acceptance. The state \mathbf{k} is a departure state where the departing call leaves k_i class i calls behind, $1 \leq i \leq J$. For states of the form \mathbf{k} where no decision needs to be made, we set $A(\mathbf{k}) = \{0\}$.

To complete the specification of the SMDP, we need to provide transition probabilities, mean sojourn times, rewards, and costs for each state–action pair. Let $\tau(\mathbf{i}, a)$ denote the average time spent in state \mathbf{i} (until the next decision epoch) if action $a \in A(\mathbf{i})$ is chosen. Then

$$\tau(\mathbf{k}, 0) = \left(\sum_{i=1}^J \lambda_i + k_i \mu_i \right)^{-1}$$

TABLE I
PARAMETERS

Class i	μ_i	ν_i	p_i	p_i^{cell}	p_i^{call}
1	0.1	6.0	0.025	10^{-9}	10^{-2}
2	1.0	1.5	0.100	10^{-9}	10^{-2}

and

$$\tau((\mathbf{k}, j), a) = \begin{cases} \left(\sum_{i=1}^J \lambda_i + k_i \mu_i \right)^{-1}, & a = 0 \\ \left(\sum_{i=1}^J \lambda_i + k_i \mu_i + \mu_j \right)^{-1}, & a = 1. \end{cases}$$

Let $p(\mathbf{i}, \mathbf{i}', a)$ denote the transition probability from state \mathbf{i} to \mathbf{i}' if action a is chosen. Let

$$\bar{p}(\mathbf{k}, \mathbf{i}) = \begin{cases} k_l \mu_l / \left(\sum_{i=1}^J \lambda_i + k_i \mu_i \right) & \text{if } \mathbf{i} = (\mathbf{k} - \mathbf{e}_l) \\ \lambda_l / \left(\sum_{i=1}^J \lambda_i + k_i \mu_i \right) & \text{if } \mathbf{i} = (\mathbf{k}, l) \end{cases}$$

and 0 otherwise. Then

$$p(\mathbf{k}, \mathbf{i}, 0) = \bar{p}(\mathbf{k}, \mathbf{i})$$

and

$$p((\mathbf{k}, j), \mathbf{i}, a) = \begin{cases} \bar{p}(\mathbf{k}, \mathbf{i}), & a = 0, \\ \bar{p}(\mathbf{k} + \mathbf{e}_j, \mathbf{i}), & a = 1. \end{cases}$$

2) The Conservative Approach to the Cell Loss Constraint:

We formulate a conservative approach to the cell loss constraint as follows. Let

$$C = \{ \mathbf{k} : b_i(\mathbf{k}) \leq p_i^{\text{cell}} k_i p_i \nu_i, \quad 1 \leq i \leq J \}. \quad (8)$$

Then C is the set of \mathbf{k} such that when there are k_i class i calls in the system *forever*, $1 \leq i \leq J$, the long run average cell loss rate constraints are always satisfied. The set C uniquely determines a state space I : for any j , $(\mathbf{k}, j) \in I$ if and only if $\mathbf{k} \in C$; and $\mathbf{k} \in I$ if and only if there exists a j such that $\mathbf{k} + \mathbf{e}_j \in C$. The action sets for the states on the boundary of I require minor modification: given $\mathbf{k} \in C$, if $\mathbf{k} + \mathbf{e}_j \notin C$, then $A((\mathbf{k}, j)) = \{0\}$, $1 \leq j \leq J$. In this manner, C uniquely determines an SMDP. Note that given p_j^{cell} , $1 \leq j \leq J$, and R , the set C is always finite. The optimal policy obtained by solving this SMDP is conservative in terms of cell loss constraints because it will never go into any state for *any period of time* where the cell loss constraints will be violated if we stay there *forever*. This SMDP can be solved using either value iteration or LP. We used LP. For the 2 class system whose parameters are given in Table I (with $R = 45$), using CPLEX it took 88 sec on a Sparc Server 3000 to solve the LP.

3) An Aggressive Approach to the Cell Loss Constraint:

We now describe an “aggressive” approach to the cell loss constraint. In order to be able to apply standard numerical solution procedures to find the optimal policy, we need to make the state space of the SMDP finite. We achieve this by placing an *a priori* bound on the number of calls that can

be accepted into the system. Let M be a fixed large positive integer. We consider a state space $I(M)$ of the form

$$I(M) = \left\{ (\mathbf{k}, j) : k_i \in \mathbf{Z}^+, \quad i = 1, \dots, J; \right. \\ \left. 1 \leq j \leq J; \text{ and } \sum_{i=1}^J k_i \leq M \right\} \\ \cup \left\{ \mathbf{k} : k_i \in \mathbf{Z}^+, \quad i = 1, \dots, J \text{ and } \sum_{i=1}^J k_i \leq M \right\}.$$

For states (\mathbf{k}, j) such that $\sum_{i=1}^J k_i = M$, we let $A((\mathbf{k}, j)) = \{0\}$.

Since M was chosen arbitrarily to make the problem finite, we need to solve a series of problems with increasing M 's, until the associated optimal policy and reward stop changing. Then the optimal policy for $M = \infty$ will have been obtained.

Let $r(\mathbf{k}, j, a)$ denote the reward earned in state (\mathbf{k}, j) when action a is chosen. (The reward earned in states of the form \mathbf{k} is 0.) Then

$$r(\mathbf{k}, j, a) = \begin{cases} 0, & a = 0 \\ w_j, & a = 1. \end{cases}$$

We also define “costs” which will play a role in the constraint on loss probability.

Let

$$\bar{l}_\ell(\mathbf{k}) = b_\ell(\mathbf{k}) \tau((\mathbf{k}), 0), \quad \ell = 1, \dots, J.$$

We define

$$l_\ell((\mathbf{k}), 0) = \bar{l}_\ell(\mathbf{k})$$

and

$$l_\ell((\mathbf{k}, j), a) = \begin{cases} \bar{l}_\ell(\mathbf{k}), & a = 0 \\ \bar{l}_\ell(\mathbf{k} + \mathbf{e}_j), & a = 1, j = 1, \dots, J. \end{cases}$$

As defined above, $l_\ell(\mathbf{i}, a)$ is the expected amount of class ℓ traffic lost until the next decision epoch. Let

$$\bar{g}_\ell(\mathbf{k}) = p_\ell k_\ell \nu_\ell \tau((\mathbf{k}), 0) \\ g_\ell(\mathbf{k}, 0) = \bar{g}_\ell(\mathbf{k})$$

and

$$g_\ell((\mathbf{k}, j), a) = \begin{cases} \bar{g}_\ell(\mathbf{k}), & a = 0 \\ \bar{g}_\ell(\mathbf{k} + \mathbf{e}_j), & a = 1, j = 1, \dots, J. \end{cases}$$

Then $g_\ell(\mathbf{i}, a)$ is the expected amount of class ℓ traffic to arrive until the next decision epoch.

We can again use the results of [7], which enable us to obtain the optimal control from the following linear program:

$$\text{maximize} \quad \sum_{\mathbf{i} \in I(M)} \sum_{a \in A(\mathbf{i})} r(\mathbf{i}, a) z_{\mathbf{i}a} \quad (9)$$

subject to

$$\sum_{a \in A(\mathbf{j})} z_{\mathbf{j}a} - \sum_{\mathbf{i} \in I(M)} \sum_{a \in A(\mathbf{i})} p(\mathbf{i}, \mathbf{j}, a) z_{\mathbf{i}a} = 0, \quad \mathbf{j} \in I(M) \quad (10)$$

$$\sum_{\mathbf{i} \in I(M)} \sum_{a \in A(\mathbf{i})} [l_\ell(\mathbf{i}, a) - p_\ell^{\text{cell}} g_\ell(\mathbf{i}, a)] z_{\mathbf{i}a} \leq 0, \quad \ell = 1, \dots, J \quad (11)$$

$$\sum_{\mathbf{i} \in I(M)} \sum_{a \in A(\mathbf{i})} \tau(\mathbf{i}, a) z_{\mathbf{i}a} = 1 \quad (12)$$

$$z_{\mathbf{i}a} \geq 0, \quad \mathbf{i} \in I(M), a \in A(\mathbf{i}). \quad (13)$$

As was the case in Section II-B6, the above LP is clearly feasible, because $z_{\mathbf{0},0} = [\tau(\mathbf{0},0)]^{-1} = \sum_{i=1}^J \lambda_i$, and $z_{\mathbf{i}a} = 0$ otherwise is a feasible solution. Let $\phi(\mathbf{k}, j)$ denote the probability that action 1 (accept) is chosen in state (\mathbf{k}, j) . Given an optimal solution, z , of the LP, we obtain an optimal randomized stationary policy as

$$\phi(\mathbf{k}, j) = \frac{z(\mathbf{k}, j), 1}{z(\mathbf{k}, j), 0 + z(\mathbf{k}, j), 1}$$

if $z(\mathbf{k}, j), 0 + z(\mathbf{k}, j), 1 > 0$, and $\phi(\mathbf{k}, j) = 0$ otherwise. Again, using the parameters in Table I (with $R = 45$) and using CPLEX on a Sparc Server 3000, it took 402 sec to solve the LP.

4) *Cell Level Control*: Both the aggressive and conservative approaches involve J cell level constraints, one for each class. In both cases, it is sometimes possible to transform the problem to one involving one constraint using cell level control, which consists of deciding how many cells of each class are lost. This transformation will result in improved performance of the associated optimal control. For $J = 2$, (7) presents the cell loss rate under a “no priority” assumption. Not surprisingly, typically only one of the two cell level constraints is tight in the solution. Thus, by giving priority to the class whose constraint is tight, it seems clear that we can improve the solution. Related results and discussion are contained in [3].

Consider the conservative approach for $J = 2$, for which the “acceptance region,” C , is given by (8). Let

$$b(\mathbf{k}) = \sum_{n_1=0}^{k_1} \sum_{n_2=0}^{k_2} \psi(\mathbf{k}, \mathbf{n}) [n_1 \nu_1 + n_2 \nu_2 - R]^+. \quad (14)$$

There is a cell level control (it will depend on k_1, k_2) that can achieve

$$b_i(\mathbf{k}) \leq p_i^{\text{cell}} k_i p_i \nu_i, \quad i = 1, 2 \quad (15)$$

if

$$b(k_1, k_2) \leq p_1^{\text{cell}} k_1 p_1 \nu_1 + p_2^{\text{cell}} k_2 p_2 \nu_2 \quad (16a)$$

$$b(k_1, 0) \leq p_1^{\text{cell}} k_1 p_1 \nu_1 \quad (16b)$$

and

$$b(0, k_2) \leq p_2^{\text{cell}} k_2 p_2 \nu_2. \quad (16c)$$

Furthermore, if $p_1^{\text{cell}} = p_2^{\text{cell}}$, then (16a) implies (16b) and (16c), and the two constraints of (15) are replaced by the one constraint (16a). In this case, the acceptance region becomes

$$C = \{(k_1, k_2) : b(\mathbf{k}) \leq p^{\text{cell}} [k_1 p_1 \nu_1 + k_2 p_2 \nu_2]\} \quad (17)$$

where $p^{\text{cell}} = p_1^{\text{cell}} = p_2^{\text{cell}}$.

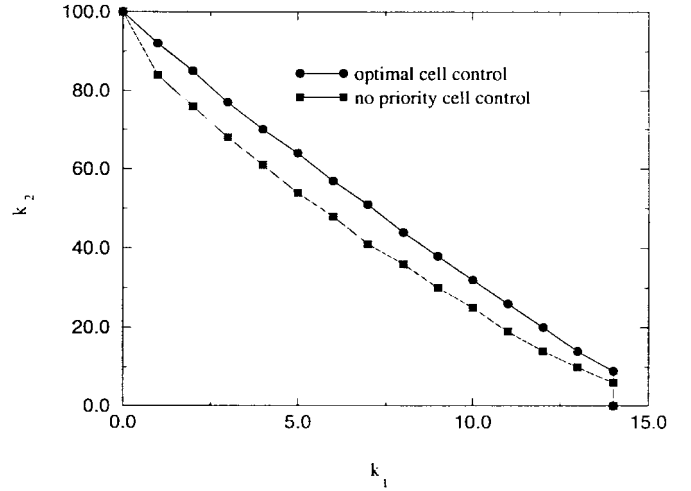


Fig. 2. Acceptance regions ($p_1^{\text{cell}} = p_2^{\text{cell}} = 10^{-9}$).

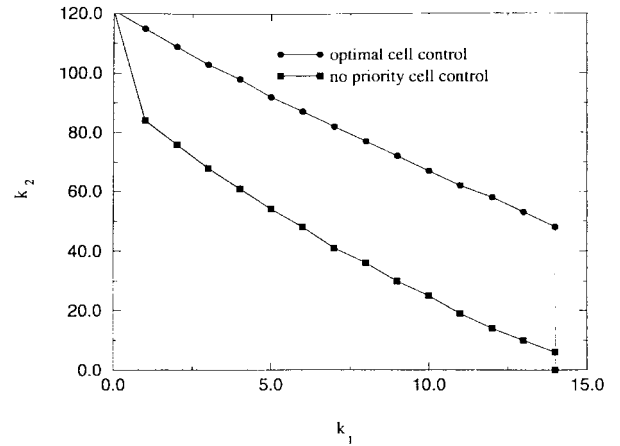


Fig. 3. Acceptance region ($p_1^{\text{cell}} = 10^{-9}, p_2^{\text{cell}} = 10^{-7}$).

The aggressive case is a bit more complicated because the constraints involve the stationary distribution over all states. Consider the constraint

$$\sum_{\mathbf{i} \in I} \sum_{a \in A(\mathbf{i})} [l_1(\mathbf{i}, a) + c_2(\mathbf{i}, a) - p_1^{\text{cell}} g_1(\mathbf{i}, a) - p_2^{\text{cell}} g_2(\mathbf{i}, a)] z_{\mathbf{i}a} \leq 0. \quad (18)$$

For the case $p_1^{\text{cell}} = p_2^{\text{cell}}$, if we solve the LP (9), (10), (12), (13), (18), we can find a cell level control such that the constraints (11) are satisfied. A related result holds with $p_1^{\text{cell}} \neq p_2^{\text{cell}}$.

An indication of the advantage of cell level control can be seen by comparing the acceptance regions for the conservative scheme. Fig. 2 displays the acceptance regions associated with one and two constraints, as given by (17) and (18), for the parameter values of the two classes given in Table I (with $R = 45$).

In Fig. 3, where we use $p_1^{\text{cell}} = 10^{-9}$ and $p_2^{\text{cell}} = 10^{-7}$, the advantage of cell level control is more dramatic.

IV. FEASIBILITY OF CALL ADMISSION CONTROLS

In this section we discuss the feasibility problem which focuses not only on cell level QoS but also on call level QoS. Given parameters $R, \lambda_1, \lambda_2, \mu_1, \mu_2, p_1, p_2, \nu_1, \nu_2$, and QoS constraints $p_1^{\text{cell}}, p_2^{\text{cell}}, p_1^{\text{call}}, p_2^{\text{call}}$, does there exist a CAC under which the QoS constraints are satisfied? The exact version of this problem was shown, in Section II-B7, to reduce to feasibility of an LP. Using the NCD results of Section III, we can also obtain an LP for the reduced problem. For the aggressive approach, the development is very similar to that in Section II-B-7, so we can be brief.

We define the set of arrival states corresponding to class j calls as

$$I_j^A(M) = \left\{ (\mathbf{k}, j) : k_i \in \mathbb{Z}^+, 1 \leq i \leq J; \sum_{i=1}^J k_i \leq M \right\}, \\ 1 \leq j \leq J.$$

The resulting LP is

$$\text{maximize} \quad \sum_{\mathbf{i} \in I(M)} \sum_{a \in A(\mathbf{i})} r(\mathbf{i}, a) z_{\mathbf{i}a}$$

subject to

$$\sum_{a \in A(\mathbf{j})} z_{\mathbf{j}a} - \sum_{\mathbf{i} \in I(M)} \sum_{a \in A(\mathbf{i})} p(\mathbf{i}, \mathbf{j}, a) z_{\mathbf{i}a} = 0, \quad \mathbf{j} \in I(M)$$

$$\sum_{\mathbf{i} \in I(M)} \sum_{a \in A(\mathbf{i})} [\ell_\ell(\mathbf{i}, a) - p_\ell^{\text{cell}} g_\ell(\mathbf{i}, a)] z_{\mathbf{i}a} \leq 0, \quad \ell = 1, \dots, J$$

$$\sum_{\mathbf{i} \in I_j^A(M)} z_{\mathbf{i}0} - p_j^{\text{call}} \lambda_j \leq 0, \quad 1 \leq j \leq J$$

$$\sum_{\mathbf{i} \in I(M)} \sum_{a \in A(\mathbf{i})} \tau(\mathbf{i}, a) z_{\mathbf{i}a} = 1 \\ z_{\text{in}} \geq 0, \quad \mathbf{i} \in I(M), a \in A(\mathbf{i}).$$

The above LP will be feasible if and only if there is a CAC under which both the cell and call level QoS constraints are met. For $J = 2$, holding all parameters other than λ_1 and λ_2 fixed, we solved for the set of (λ_1, λ_2) for which the LP is feasible. We call this set the feasible region. An example of the feasible region is illustrated in Fig. 5 for the parameter values in Table I. (Other examples are contained in [20].) Remarkably, the boundary of this region exhibits near linearity. As a consequence of this (almost) linearity, the feasible region can be (approximately) determined from its corner points. These corner points can be obtained by solving two one-dimensional (single class) problems.

An LP for the conservative approach with call level QoS constraints can also be formulated and solved. The feasible region for the conservative approach is also indicated in Fig. 5. Note that the feasible region for the aggressive approach is substantially larger than for the conservative approach. We thus focus on the aggressive approach in the rest of the paper.

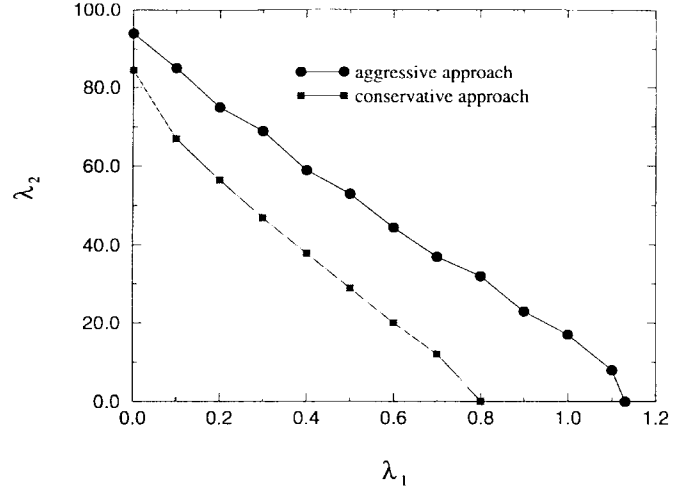


Fig. 5. Feasible regions: $p_i^{\text{cell}} = 10^{-9}$, $p_i^{\text{call}} = 0.01$, $i = 1, 2$.

V. THE SINGLE CLASS SYSTEM WITH CELL AND CALL QoS

Both cell and call level QoS are considered in the following study of a single class system. The insights gained from this study provide the basis for the techniques for multiclass systems.

A. The Model

The system and traffic model considered here is that obtained using NCD, with $J = 1$. Recall that for the system, the key parameters for the single class are denoted by λ, μ, p, ν ; the QoS parameters are p^{cell} and p^{call} ; and the link rate is R .

For the one-dimensional reduced Markov chain, given a stationary call admission control policy, various steady-state performance measures of the system can be easily calculated. In this and the following sections, we do not specify any rewards or costs to be optimized. Our goal here is feasibility and robustness, so we do not provide any cost function. If we received a reward for each accepted call, we would try to admit as many calls as possible subject to the two QoS requirements. Although the *optimal* admission control policy in this context normally has a randomized threshold due to the cell and call level QoS constraints, nonrandomized threshold policies, which are simpler, are sufficiently accurate for our purpose here.

B. Some Basic Relations

Let $\phi(\kappa)$ be an admission control policy based on a threshold κ , i.e., a new call is admitted if there are fewer than κ calls in progress, otherwise it is rejected. Based on the above model, we can calculate the following performance measures under policy $\phi(\kappa)$. The mean cell arrival rate with k calls permanently in progress is $s(k) = k\nu p$. The mean cell loss rate with k calls permanently in progress is $b(k) = \sum_{n=0}^k \binom{k}{n} p^n (1-p)^{k-n} [n\nu - R]^+$. Let $g(\kappa) = \sum_{k=0}^{\kappa} (\lambda/\mu)^k / k!$. The stationary distribution for the number of calls in progress is $P(k \text{ calls}) = [(\lambda/\mu)^k / k!] / g(\kappa)$, $k = 0, 1, \dots, \kappa$. The average cell arrival rate is $s_\kappa = [(\sum_{k=1}^{\kappa} (\lambda/\mu)^k / k!) s(k)] / g(\kappa)$. The average cell loss rate is $b_\kappa = [\sum_{k=1}^{\kappa} ((\lambda/\mu)^k / k!) b(k)] / g(\kappa)$. The average

cell loss ratio is

$$L(\lambda, \kappa, R) = \frac{b_\kappa}{s_\kappa} = \frac{\sum_{k=1}^{\kappa} \frac{(\lambda/\mu)^k}{k!} b(k)}{\sum_{k=1}^{\kappa} \frac{(\lambda/\mu)^k}{k!} s(k)}. \quad (19)$$

The call blocking probability is given by the Erlang loss formula $B(\lambda, \kappa) = [(\lambda/\mu)^\kappa / \kappa!] / g(\kappa)$.

C. Feasibility of Call Arrival Rate

Similarly to the treatment in Section IV, we can define the feasibility of the call arrival rate for a single class system as follows: a call arrival rate λ is said to be feasible if there exists an admission control policy from the class of threshold policies that meets both QoS requirements. That is, λ is feasible if there exists κ such that $\phi(\kappa)$ ensures both cell and call QoS. Recall that our cell and call QoS requirements are, respectively: $L(\lambda, \kappa, R) \leq p^{\text{cell}}$, and $B(\lambda, \kappa) \leq p^{\text{call}}$. For a given λ , let

$$\kappa_1(\lambda) = \max \{ \kappa : L(\lambda, \kappa) \leq p^{\text{cell}} \} \quad \text{and} \\ \kappa_2(\lambda) = \min \{ \kappa : B(\lambda, \kappa) \leq p^{\text{call}} \}.$$

If $\kappa_2(\lambda) \leq \kappa_1(\lambda)$, then there exists a threshold-type policy $\phi(\kappa)$ with $\kappa_2(\lambda) \leq \kappa \leq \kappa_1(\lambda)$, such that under $\phi(\kappa)$ both cell and call QoS requirements are met, and we say λ is feasible.

Given p^{cell} , p^{call} and the traffic parameters described in Subsection A, there is a maximum feasible call arrival rate, λ_{\max} , such that for any $\lambda > \lambda_{\max}$, $\kappa_2(\lambda) > \kappa_1(\lambda)$. It is not surprising that for most cases of practical interest, $\kappa_2(\lambda_{\max}) = \kappa_1(\lambda_{\max})$, which is assumed throughout this paper.

D. An Effective Bandwidth

Let $e = R/\kappa_1(\lambda_{\max})$. Then e is a measure of the resources a call requires to satisfy *both* cell level and call level QoS requirements. Note that e is independent of the call arrival rate λ .

Note that a traditional effective bandwidth definition (see [6], for instance) assumes that calls last in perpetuity and addresses only cell level QoS

$$e_{\text{static}} = \frac{R}{\kappa_{\text{static}}}, \quad \text{where} \\ \kappa_{\text{static}} = \max \left\{ k : \frac{b(k)}{s(k)} \leq p^{\text{cell}} \right\}.$$

Note that $e \leq e_{\text{static}}$. Since e_{static} ignores call level dynamics, it is more conservative.

Example 1: Consider a system with $R = 45.0$, and homogeneous sources with parameters corresponding to class 1 sources as listed in Table I, i.e., $\mu = 0.1$, $\nu = 6.0$, $p = 0.025$, $p^{\text{cell}} = 10^{-9}$, and $p^{\text{call}} = 0.01$. In Fig. 6, we illustrate the relation between $\kappa_1(\lambda)$, $\kappa_2(\lambda)$, and λ_{\max} . $\kappa_1(\lambda)$ and $\kappa_2(\lambda)$ intersect at 20 when $\lambda = 1.125$. Hence, $\lambda_{\max} = 1.125$, and $e = \frac{45}{20} = 2.25$. Note that $\kappa_{\text{static}} = 14$, hence $e_{\text{static}} = 3.21$.

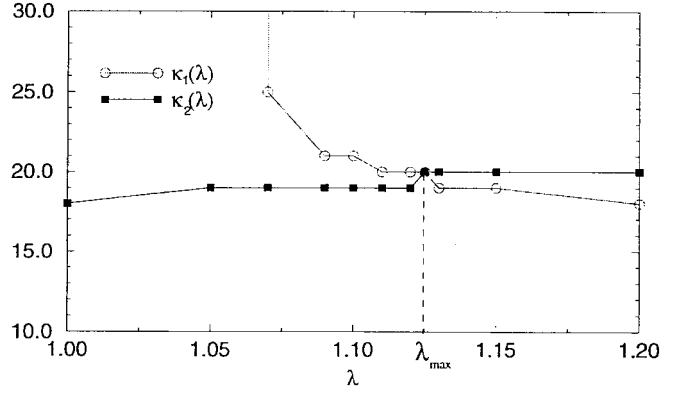


Fig. 6. Calculating effective bandwidth for class 1 sources.

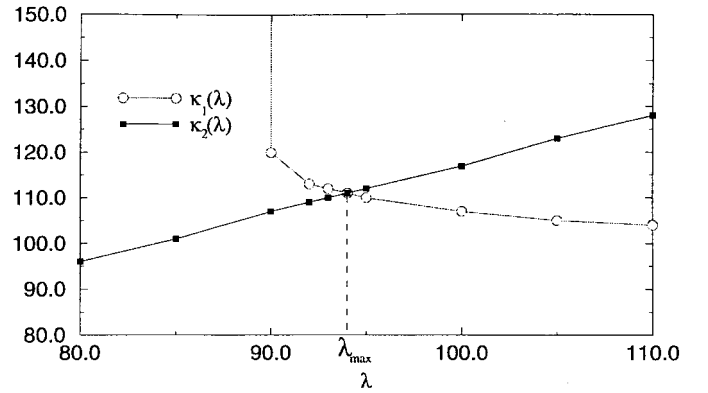


Fig. 7. Calculating effective bandwidth for class 2 sources.

Example 2: Now consider a system with the same R , p^{cell} , p^{call} as in Example 1, and sources corresponding to class 2, see Table I, i.e., $\mu = 1.0$, $\nu = 1.5$, and $p = 0.1$. The relation between $\kappa_1(\lambda)$, $\kappa_2(\lambda)$, and λ_{\max} are illustrated in Fig. 7. $\kappa_1(\lambda)$ and $\kappa_2(\lambda)$ intersect at 111 when $\lambda = 94.0$. Hence, $\lambda_{\max} = 94.0$, and $e = \frac{45}{111} = 0.4054$. Now $\kappa_{\text{static}} = 100$, and $e_{\text{static}} = 0.45$.

VI. ASYMPTOTIC ANALYSIS

An asymptotic analysis of a single class system provides fundamental insights into the joint behavior of the cell loss ratio and call blocking probability for purposes of sizing and operations. The investigation is in the asymptotic framework of large systems, i.e., as $(\lambda, \kappa, R) \rightarrow \infty$ in a manner consistent with practical QoS requirements. Specifically, p^{cell} is expected to be in the range 10^{-6} – 10^{-9} , while p^{call} is expected to be in the neighborhood of 10^{-2} . These numbers suggest the following important dichotomy: cell loss ratios decay exponentially in the large parameter, say κ , while call blocking probabilities decay polynomially, more specifically as $1/\sqrt{\kappa}$. While both elements are separately recognized in the literature (see, for instance, [24]), we do not know of any prior analysis in which both elements are simultaneously present. The loadings at the cell and call levels are required to be such that the ultimate cell and call performances are, respectively, exponential and polynomial in κ . We obtain such loading guidelines. We also

obtain the coefficients associated with the exponential and polynomial behaviors. Thus the contributions here are both qualitative and quantitative.

First, we write the expression for the cell loss ratio in (19) as $L(\lambda, \kappa, C)$, where $C = R/\nu$ by definition. Also, we assume for convenience that C is an integer, and select the unit of time to be such that $\mu = 1$. Then

$$L(\lambda, \kappa, C) = \text{Num}/\text{Den} \quad (20)$$

where

$$\begin{aligned} \text{Num} &= 0 \quad \text{if } \kappa \leq C \\ &= \sum_{k=C+1}^{\kappa} \frac{\lambda^k}{k!} \sum_{n=C+1}^k (n-C) \binom{k}{n} p^n (1-p)^{k-n} \\ &\quad \text{if } \kappa > C \end{aligned} \quad (21)$$

and

$$\text{Den} = p \sum_{k=0}^{\kappa} k \lambda^k / k!. \quad (22)$$

Note that we may assume that $\kappa > C$, as otherwise $L = 0$.

As mentioned above, we let $(\lambda, \kappa, C) \rightarrow \infty$, while p and ν are held fixed. For loading at the cell level, we assume that the system is underloaded, i.e.,

$$\rho \triangleq \frac{\lambda p}{C} < 1 \quad (23)$$

and that critical loading holds at the call level, i.e.,

$$\gamma \triangleq \left(1 - \frac{\kappa}{\lambda}\right) \sqrt{\lambda} = O(1) \quad (24)$$

i.e., γ is bounded. At the loss of some small generality, we will make the convenient assumption that γ is a fixed constant, which can be either positive or negative. Hence, $\kappa = \lambda - \gamma \sqrt{\lambda}$, so that, to leading order, $\kappa/\lambda \sim 1$.

Here we give an overview of the analysis, which includes the salient elements, with the details provided in the Appendix. We show that

$$\begin{aligned} \text{Num} &= \sum_{l=0}^{\kappa-C} \frac{l(\lambda p)^{C+l}}{(C+l)!(\kappa-C-l)!} \\ &\quad \times \int_0^{\infty} e^{-u} [\lambda(1-p) + u]^{\kappa-C-l} du \end{aligned} \quad (25)$$

where the integral is obtained by making use of Euler's representation, $n! = \int_0^{\infty} e^{-u} u^n du$, and the binomial theorem. Now

$$(C+l)! = C!(C+1) \cdots (C+l) \geq C! C^l \quad (26)$$

and also the bound is asymptotically exact as $C \rightarrow \infty$ and $l = O(1)$. After substituting (26) in (25), we show that

$$\begin{aligned} \text{Num} &\sim \frac{1}{C!} \sum_{l=0}^{\kappa-C} \frac{l(\lambda p)^{C+l}}{C^l(\kappa-C-l)!} \\ &\quad \times \int_0^{\infty} e^{-u} [\lambda(1-p) + u]^{\kappa-C-l} du \\ &= \frac{(\lambda p)^C}{C!(\kappa-C)!} \left[(\kappa-C) \int_0^{\infty} \int_0^{\infty} e^{-(u+v)} \right. \\ &\quad \times \rho v [\lambda(1-p) + u + \rho v]^{\kappa-C-1} du dv \Big] \\ &= \frac{(\lambda p)^C \rho}{C!(\kappa-C)!} \frac{\partial I}{\partial \rho} \end{aligned} \quad (27)$$

where

$$I = \int_0^{\infty} \int_0^{\infty} e^{-(u+v)} [\lambda(1-p) + u + \rho v]^{\kappa-C} du dv. \quad (28)$$

This is an important integral representation.

We obtain an asymptotic expression for the integral I for the scaling in (23) and (24):

$$I \sim \frac{\{\lambda(1-p)\}^{\kappa-C}}{(1-\alpha)(1-\rho\alpha)} \quad (29)$$

where

$$\alpha \triangleq \frac{1-p/\rho}{1-p}. \quad (30)$$

The constant α arises since $(\kappa-C)/\{\lambda(1-p)\} \sim \alpha$ from the scaling. Note that $0 < \alpha < 1$, where the positivity is due to $\kappa > C$.

From (29)

$$\frac{\partial I}{\partial \rho} \sim \frac{\alpha \{\lambda(1-p)\}^{\kappa-C}}{(1-\alpha)(1-\rho\alpha)^2}. \quad (31)$$

Substitution in (27) together with use of Stirling's formula for $C!$ and $(\kappa-C)!$ gives the asymptotic expression for Num.

Turning to Den in (22), we find that

$$\text{Den} = \frac{p \lambda^{\kappa}}{(\kappa-1)!} \frac{1}{B(\lambda, \kappa-1)} \quad (32)$$

where $B(\lambda, \kappa-1)$ is the Erlang loss formula. Since the call level loading is critical, we know from prior results [12] that to leading order $B(\lambda, \kappa-1)/B(\lambda, \kappa) \sim 1$, and that

$$\frac{1}{B(\lambda, \kappa)} \sim \sqrt{\lambda} W_0(\gamma) \quad (33)$$

where $W_0(\gamma)$ is an $O(1)$ constant which is obtained from the normal distribution and depends only on γ . In fact

$$e^{-\gamma^2/2} W_0(\gamma) = \left(\frac{\pi}{2}\right)^{1/2} \text{erfc}\left(\frac{\gamma}{\sqrt{2}}\right)$$

where erfc is the complementary error function.

Finally, substituting the asymptotic expressions for Num and Den, we obtain

$$L(\lambda, \kappa, C) \sim \frac{A e^{-\delta \kappa}}{\kappa^2} \quad (34)$$

where

$$A = A(\rho, p, \gamma) = \frac{1}{\sqrt{2\pi}} \frac{\rho\alpha}{p(1-\alpha)(1-\rho\alpha)^2} \\ \times \frac{1}{\sqrt{(1-p/\rho)p/\rho}} \cdot \frac{1}{W_0(\gamma)}$$

and

$$\delta = \delta(\rho, p) = \log \left[\left(\frac{1-p/\rho}{1-p} \right)^{1-p/\rho} \left(\frac{1}{\rho} \right)^{p/\rho} \right]. \quad (35)$$

It is easy to verify that for all (ρ, p) such that $0 < \rho < 1$ and $0 < p < 1$, $\delta(\rho, p) > 0$. Clearly, $A(\rho, p, \gamma) = O(1)$. These facts prove the important result that for the asymptotic scaling in (23) and (24), the cell loss ratio $L(\lambda, \kappa, C)$ is exponentially small in the large parameter κ , with δ the constant in the exponential. The asymptotic call blocking probability for our scaling has already been obtained in (33), in which the $1/\sqrt{\lambda}$ (and therefore $1/\sqrt{\kappa}$) behavior is exhibited, and $1/W_0(\gamma)$ is the relevant constant.

VII. CAC DESIGN

We now return to the two-class problem. We seek a CAC that satisfies all QoS requirements. For feasible parameters, such a CAC can be obtained from the LP. We are seeking, however, a simple and robust CAC.

We consider all parameters except for λ_1 and λ_2 to be fixed, and consider a point (λ_1, λ_2) inside the feasible region described in Section IV and depicted in Fig. 5. We consider situations where class 1 is “burstier.”

The two-class system that we study has $R = 45$ and other parameters specified in Table I. We know that when (λ_1, λ_2) is on or very close to the boundary of the feasible region, the call admission policies that meet all the QoS requirements necessarily have complicated structures. For just this reason, the engineered operating points are designed not to be too close to the boundary. Under the assumption that (λ_1, λ_2) are not very close to the boundary, the question of interest in this paper is whether we can find a *simple* connection admission policy that meets *both cell and call QoS requirements*. We want our policy to be robust in the sense that if the operating point drifts away from the engineered loads due to the unexpected high arrival rate of one class, the admission policy should be able to protect the other class.

The most widely studied call admission policy type is CS. Although it is easy to implement, it always favors calls requiring less bandwidth capacity, thus it can drive the blocking probability of calls with a larger bandwidth requirement up when the arrival rates of the calls with smaller bandwidth requirement are high. A very undesirable situation is when the class 1 blocking probability exceeds the QoS constraint, while the blocking probability for class 2 calls is much lower than its constraint. With CS, regardless which class exceeds its engineered load, class 1 will suffer in terms of call blocking probability.

Trunk reservation (TR) policies have been known to be able to provide protection to wideband calls. With properly chosen trunk reservation parameters, class 1 calls will not have unacceptably high call blocking probabilities due to high arrival

rates of class 2 calls. However, the traditional trunk reservation gives priority to a fixed class, and hence it cannot protect the other class when the favored class has high arrival rate. Of particular interest is the following trunk reservation policy which we study numerically in Section IX: for single link systems with two classes of calls, narrowband and wideband, if we admit a new call (regardless of its class) only when the spare capacity in the system is at least the bandwidth of wideband calls, then this special trunk reservation policy balances the call blocking probabilities of the two classes, which can be easily seen with PASTA [16]. Furthermore, this policy is optimal among all the trunk reservation policies that balance the call blocking probabilities in the sense that the call blocking probabilities are the smallest. However, because this policy balances the call blocking probabilities, if one class has high call arrival rate, both classes will have high call blocking probabilities. In other words, it is not robust, which is far from desirable in many situations.

It is intuitively apparent that if (λ_1, λ_2) are small enough (close to the origin), a complete sharing policy will be sufficient. The interesting case is when (λ_1, λ_2) is not very close to the boundary of the feasible region and also not small. In this case, we believe a simple policy based on virtual partitioning with properly chosen parameters gives satisfactory performance.

VIII. VIRTUAL PARTITIONING

We now describe a call admission policy based on the notion of VP. Let e_1 and e_2 be the bandwidth requirements of class 1 and 2 calls, and K_1 and K_2 be the partitioning parameters which are two positive integers such that $K_1 e_1 + K_2 e_2 \geq R$. A call admission policy based on VP is summarized as follows. When a call of class 1 arrives and finds (k_1, k_2) calls in progress, it is accepted if

$$\begin{aligned} k_1 < K_1 \quad \text{and} \quad e_1 k_1 + e_2 k_2 \leq R - e_1, \quad \text{or} \\ k_1 \geq K_1 \quad \text{and} \quad e_1 k_1 + e_2 k_2 \leq R - t_2 e_2 - e_1 \end{aligned}$$

where $t_2 e_2$ is the bandwidth reserved for (underloaded) class 2 calls. Similarly, a class 2 call is accepted if

$$\begin{aligned} k_2 < K_2 \quad \text{and} \quad e_1 k_1 + e_2 k_2 \leq R - e_2, \quad \text{or} \\ k_2 \geq K_2 \quad \text{and} \quad e_1 k_1 + e_2 k_2 \leq R - t_1 e_1 - e_2 \end{aligned}$$

where $t_1 e_1$ is the bandwidth reserved for (underloaded) class 1 calls.

Note that the call admission is performed on a set in (k_1, k_2) space defined by $k_1 e_1 + k_2 e_2 \leq R$. The motivation for selecting this set is derived from the linearity implicit in the notion of effective bandwidths. Admittedly there is no sound theoretical basis for the linearity at this time. In the absence of such a theory, we take the precautionary step of verifying in our numerical investigations that cell level QoS requirements are satisfied by our admission control policies.

By choosing parameters K_1 and K_2 , we partition the bandwidth between the two classes. The trunk reservation parameters t_1 and t_2 allow us to block calls from the overloaded class so as to reserve bandwidth for the underloaded class. The

TABLE III
NUMERICAL RESULTS FOR CASE 1

Call arrival rates (λ_1, λ_2)	Call blocking Probabilities (B_1, B_2)		
	CS	TR	VP
(0.60, 35.0)	(0.0064, 0.0008)	(0.0031, 0.0031)	(0.0049, 0.0033)
(0.72, 35.0)	(0.0184, 0.0026)	(0.0096, 0.0096)	(0.0153, 0.0076)
(0.60, 45.0)	(0.0232, 0.0032)	(0.0106, 0.0106)	(0.0083, 0.0202)
(0.66, 38.5)	(0.0174, 0.0024)	(0.0085, 0.0085)	(0.0110, 0.0110)
(0.55, 33.0)	(0.0026, 0.0003)	(0.0013, 0.0013)	(0.0021, 0.0012)

nature of the policy is to give the underloaded class higher priority, and the consequence is that the underloaded class is protected.

We expect t_1 and t_2 to be small nonnegative integers. When $t_1 = t_2 = 0$, the policy becomes a complete sharing policy over the whole admissible region (regardless of K_1 and K_2). If in this case the call blocking probabilities for *both* classes are still greater than the allowed limits (say, 1%), then the arrival rates are too high for there to exist any feasible policy. When t_1 and t_2 are not both zero, the policy is complete sharing on the set $\{(k_1, k_2) : k_1 \leq K_1, k_2 \leq K_2\}$, and dynamically prioritized outside the set. With complete sharing, the less bursty traffic enjoys lower call blocking, and consequently when the set in which complete sharing applies is bigger, fewer calls of the less bursty traffic are blocked.

IX. NUMERICAL RESULTS FOR VP

In our numerical study, we attempt to show that we can use VP to design simple call admission policies such that when the system is subject to load at or below the engineered call arrival rates (λ_1^e, λ_2^e), the admission policy will meet all the cell and call QoS requirements; and when the system is subject to load above the engineered loads due to a high arrival rate from one class, the other class is protected. More specifically, we observe the call blocking probabilities of the two classes, denoted by (B_1, B_2) , under CS, TR, and VP to illustrate the robustness of VP. The following load scenarios are studied: both classes are below the engineered load; both classes are at the engineered load; one of the two classes is at the engineered load while the other is at least 20% higher; both classes are 10% higher than the engineered loads.

The system and traffic sources are the ones specified in Section VII. From the analysis of the two single class problems in Section V-D, we have $c_1 = 2.25$ and $c_2 = 0.4054$, which are used for all the cases in this section.

A. Case 1

The first set of numerical results are for the engineered load: $(\lambda_1^e, \lambda_2^e) = (0.6, 35.0)$. The bandwidth partitioning parameters used are $(K_1, K_2) = (9, 42)$, and the trunk reservation parameters used are $(t_1, t_2) = (2, 0)$. The call blocking probabilities for several different loads are listed in Table III.

When the load is below or at the engineered load, all three policies give satisfactory call blocking probabilities. However, when one of the two classes has load higher than the engineered load, class 1 always suffers under CS and the blocking probabilities are driven above the allowed limit 1%.

TABLE IV
NUMERICAL RESULTS FOR CASE 2

Call arrival rates (λ_1, λ_2)	Call blocking Probabilities (B_1, B_2)		
	CS	TR	VP
(0.20, 70.0)	(0.0150, 0.0018)	(0.0048, 0.0048)	(0.0004, 0.0059)
(0.24, 70.0)	(0.0241, 0.0031)	(0.0079, 0.0079)	(0.0010, 0.0096)
(0.20, 84.0)	(0.1049, 0.0153)	(0.0327, 0.0327)	(0.0018, 0.0386)
(0.22, 77.0)	(0.0536, 0.0073)	(0.0171, 0.0171)	(0.0013, 0.0205)
(0.18, 63.0)	(0.0026, 0.0003)	(0.0009, 0.0009)	(0.0001, 0.0011)

TABLE V
NUMERICAL RESULTS FOR CASE 3

Call arrival rates (λ_1, λ_2)	Call blocking Probabilities (B_1, B_2)		
	CS	TR	VP
(0.9, 10.0)	(0.0038, 0.0005)	(0.0030, 0.0030)	(0.0038, 0.0005)
(1.08, 10.0)	(0.0156, 0.0022)	(0.0128, 0.0128)	(0.0156, 0.0022)
(0.90, 15.0)	(0.0073, 0.0010)	(0.0052, 0.0052)	(0.0072, 0.0017)
(0.99, 11.0)	(0.0093, 0.0013)	(0.0073, 0.0073)	(0.0093, 0.0014)
(0.81, 9.0)	(0.0013, 0.0002)	(0.0010, 0.0010)	(0.0013, 0.0002)

TR, in general, gives better performance than CS. However, when $(\lambda_1, \lambda_2) = (0.6, 45.0)$, TR gives balanced call blocking probabilities of 0.0106. Hence, class 1 suffers because class 2 has an arrival rate higher than the engineered load. VP is able to protect the underloaded class. Because the engineered load is near the middle of the boundary of the feasible region, we need to protect each class when the other class exceeds its engineered load. Since on the complete sharing set class 2 is favored, by choosing the set large enough for VP we can protect class 2. Protection for class 1 in VP is achieved by choosing the complete sharing set not too large and selecting proper trunk reservation parameter t_1 , which is set to 2 in this case.

B. Case 2

The second set of numerical results are for the engineered load $(\lambda_1^e, \lambda_2^e) = (0.2, 70.0)$. The bandwidth partitioning parameters used are $(K_1, K_2) = (7, 60)$, and the trunk reservation parameters used are $(t_1, t_2) = (1, 0)$. Table IV summarizes the call blocking probabilities corresponding to the five scenarios.

Now the engineered load is near the upper-left corner of the feasible region. Again, at $(\lambda_1, \lambda_2) = (0.20, 84.0)$, class 1 suffers under CS and TR when class 2 exceeds the engineered load, while VP is able to protect it. When the class 1 arrival rate exceeds the engineered load by 20%, the impact on class 2 is very small. On the other hand, increasing the arrival rate of class 2 has more dramatic impact on class 1 traffic, hence it seems now the focus should be on protecting class 1 against class 2.

C. Case 3

The third set of results are for the engineered load $(\lambda_1^e, \lambda_2^e) = (0.9, 10.0)$. The bandwidth partitioning parameters used are $(K_1, K_2) = (10, 25)$, and the trunk reservation parameters used are $(t_1, t_2) = (3, 0)$. Table V summarizes the call blocking probabilities for the five scenarios.

This is a case where traffic is light so that VP is almost identical to CS. At $(\lambda_1, \lambda_2) = (1.08, 10.0)$, class 2 suffers under TR, while both CS and VP are able to protect class 2.

X. NEW DIRECTIONS

Although our results have been presented for a two-class system, we believe that they would apply for more than two classes. The key to our approach is the near linearity of the feasible region depicted in Fig. 5. If the feasible region is linear for the multiclass system, then the $J(\geq 2)$ class system can be analyzed using J single class systems, so that the CAC introduced here can be used for $J > 2$ as well. Further work on this topic is necessary.

In this paper we considered two versions of the cell loss constraint: conservative and aggressive. One might argue that the conservative approach is too conservative, and the aggressive approach is too aggressive, so that it would be nice to have something in between the two. An approach based on expected cell loss during the lifetime of a call in the system is considered in [19]. It is shown there that in the single-class system, the threshold based on this new approach indeed lies between the other two.

APPENDIX

Here we give some details of the asymptotic analysis in Section VI. We do three things in this Appendix. First, we prove (25); second, we make use of (26) to prove (27); and finally, we establish (29).

To prove (25) for $\kappa > C$, we note that from the expression in (21)

$$\text{Num} = \sum_{n=C+1}^{\kappa} \frac{(\lambda p)^n}{n!} (n-C) \sum_{k=0}^{\kappa-n} \frac{\{\lambda(1-p)\}^k}{k!}. \quad (\text{A1})$$

Now, we may use Euler's representation for $n!$ to obtain

$$\begin{aligned} & \sum_{k=0}^{\kappa-n} \frac{\{\lambda(1-p)\}^k}{k!} \\ &= \int_0^\infty e^{-u} \sum_{k=0}^{\kappa-n} \frac{\{\lambda(1-p)\}^k}{k!} \frac{u^{\kappa-n-k}}{(\kappa-n-k)!} du \\ &= \frac{1}{(\kappa-n)!} \int_0^\infty e^{-u} \sum_{k=0}^{\kappa-n} \binom{\kappa-n}{k} \{\lambda(1-p)\}^k u^{\kappa-n-k} du \\ &= \frac{1}{(\kappa-n)!} \int_0^\infty e^{-u} [\lambda(1-p) + u]^{\kappa-n} du. \end{aligned} \quad (\text{A2})$$

Substituting (A2) in (A1),

$$\begin{aligned} \text{Num} &= \sum_{n=C+1}^{\kappa} \frac{(\lambda p)^n}{n!} \frac{(n-C)}{(\kappa-n)!} \\ &\quad \times \int_0^\infty e^{-u} [\lambda(1-p) + u]^{\kappa-n} du \\ &= \sum_{l=0}^{\kappa-C} \frac{l(\lambda p)^{C+l}}{(C+l)!(\kappa-C-l)!} \\ &\quad \times \int_0^\infty e^{-u} [\lambda(1-p) + u]^{\kappa-C-l} du \end{aligned} \quad (\text{A3})$$

which is (25).

To prove (27), we make use of (26) in (A3) to obtain

$$\begin{aligned} \text{Num} &\sim \frac{1}{C!} \sum_{l=0}^{\kappa-C} \frac{l(\lambda p)^{C+l}}{C^l(\kappa-C-l)!} \\ &\quad \times \int_0^\infty e^{-u} [\lambda(1-p) + u]^{\kappa-C-l} du \\ &= \frac{(\lambda p)^C}{C!} \int_0^\infty e^{-u} \left[\sum_{l=0}^{\kappa-C} \frac{l \rho^l [\lambda(1-p) + u]^{\kappa-C-l}}{(\kappa-C-l)!} \right] du \end{aligned} \quad (\text{A4})$$

where the definition of ρ in (23) has been used. Now introduce Euler's representation of $l!$

$$\begin{aligned} \text{Num} &\sim \frac{(\lambda p)^C}{C!(\kappa-C)!} \int_0^\infty e^{-u} \int_0^\infty e^{-v} \sum_{l=0}^{\kappa-C} l \binom{\kappa-C}{l} \\ &\quad \times (\rho v)^l [\lambda(1-p) + u]^{\kappa-C-l} dv du \\ &= \frac{(\lambda p)^C}{C!(\kappa-C)!} \cdot \left[(\kappa-C) \int_0^\infty \int_0^\infty e^{-(u+v)} \rho v \right. \\ &\quad \times \sum_{l=1}^{\kappa-C} \binom{\kappa-C-1}{l-1} (\rho v)^{l-1} [\lambda(1-p) + u]^{\kappa-C-l} dv du \left. \right] \\ &= \frac{(\lambda p)^C \rho}{C!(\kappa-C)!} \left[(\kappa-C) \int_0^\infty \int_0^\infty e^{-(u+v)} \right. \\ &\quad \times v [\lambda(1-p) + u + \rho v]^{\kappa-C-1} dv du \left. \right]. \end{aligned} \quad (\text{A5})$$

The expression in outer brackets is observed to be $\partial I / \partial \rho$, where I is defined in (28). Hence (27) is proven.

To prove (29), note that from the definition of I in (28)

$$I = \{\lambda(1-p)\}^{\kappa-C} \int_0^\infty \int_0^\infty e^{-(u+v)} H(u, v; \kappa) du dv \quad (\text{A6})$$

where we let κ be the surrogate for the large parameters (λ, κ, C) , and

$$H(u, v; \kappa) \triangleq \exp \left[(\kappa-C) \log \left\{ 1 + \frac{u}{\lambda(1-p)} + \frac{\rho v}{\lambda(1-p)} \right\} \right]. \quad (\text{A7})$$

Now

$$\begin{aligned} H(u, v; \kappa) &= \left[\exp \left\{ \frac{(\kappa-C)u}{\lambda(1-p)} + \frac{(\kappa-C)\rho v}{\lambda(1-p)} \right\} \right] \\ &\quad \times \left[1 + O \left(\frac{(u+\rho v)^2}{\kappa} \right) \right]. \end{aligned} \quad (\text{A8})$$

Hence, noting that $(\kappa-C)/\{\lambda(1-p)\} \sim \alpha$ from the scaling

$$\begin{aligned} & \int_0^\infty \int_0^\infty e^{-(u+v)} H(u, v; \kappa) du dv \\ &= \int_0^\infty \int_0^\infty e^{-(u+v)} e^{\alpha u + \rho \alpha v} \left[1 + O \left(\frac{(u+\rho v)^2}{\kappa} \right) \right] du dv \\ &= \frac{1}{1-\alpha} \cdot \frac{1}{1-\rho \alpha} \left[1 + O \left(\frac{1}{\kappa} \right) \right]. \end{aligned} \quad (\text{A9})$$

The above together with (A6) gives (29).

ACKNOWLEDGMENT

The authors are grateful to our colleagues J. Morrison and K. G. Ramakrishnan for their invaluable help with the analysis in Section VI and the solution of linear programs and related numerical problems, respectively.

REFERENCES

- [1] J. M. Akinpelu, "The overload performance of engineered networks with nonhierarchical and hierarchical routing," *AT&T Bell Labs. Tech. J.*, vol. 63, pp. 1261–1281, 1984.
- [2] E. Altman and V. A. Gaitsogory, "Stability and singular perturbations in constrained Markov decision problems," *IEEE Trans. Automat. Contr.*, vol. 38, pp. 971–975, 1993.
- [3] N. G. Bean, "Effective bandwidths with different quality of service requirements," in *Integrated Broadband Communication Networks and Services*, V. B. Iverson, Ed. IFIP, 1993.
- [4] M. Beshai, R. Kositpaiboon, and J. Yan, "Interaction of call blocking and cell loss in an atm network," *IEEE J. Select. Areas Commun.*, vol. 12, pp. 1051–1058, 1994.
- [5] S. Borst and D. Mitra, "Virtual partitioning for robust resource sharing: Computational techniques for heterogeneous traffic," this issue, pp. xx–xx.
- [6] A. I. Elwalid and D. Mitra, "Effective bandwidth of general Markovian traffic sources and admission control of high speed networks," *IEEE/ACM Trans. Networking*, vol. 1, pp. 329–343, 1993.
- [7] E. A. Feinberg, "Constrained semi-Markov decision processes with average rewards," *ZOR-Mathematical Methods Operations Res.*, vol. 39, pp. 257–288, 1994.
- [8] R. J. Gibbens, F. P. Kelly, and P. B. Key, "A decision-theoretic approach to call admission control in ATM networks," *IEEE J. Select. Areas Commun.*, vol. 13, pp. 1101–1114, 1995.
- [9] F. Hubner and P. Tran-Gia, "Quasi-stationary analysis of a finite capacity asynchronous multiplexer with modulated deterministic input," presented at *ITC-13*, Copenhagen, 1991.
- [10] J. Y. Hui, "Resource allocation for broadband networks," *IEEE J. Select. Areas Commun.*, vol. 6, 1988.
- [11] J. M. Hyman, A. A. Lazar, and G. Pacifici, "A separation principle between scheduling and admission control for broadband switching," *IEEE J. Select. Areas Commun.*, vol. 11, pp. 605–616, 1993.
- [12] D. L. Jagerman, "Some properties of the Erlang loss function," *Bell Syst. Tech. J.*, vol. 53, pp. 525–551, 1974.
- [13] Special issue on "Advances in the fundamentals of networking—Part I," *IEEE J. Select. Area Commun.*, vol. 13, 1995.
- [14] F. P. Kelly, "Effective bandwidths at multi-class queues," *Queueing Syst.*, vol. 9, pp. 5–16, 1991.
- [15] P. Key, "Optimal control and trunk reservation in loss networks," *Probability Eng. Inform. Sci.*, vol. 4, pp. 203–242, 1990.
- [16] K. Lindberger, "Dimensioning and designing methods for integrated ATM networks," in *Proc. ITC-14*, 1994, pp. 897–906.
- [17] D. Mitra, M. I. Reiman, and J. Wang, "Robust admission control for heterogeneous ATM systems with both cell and call QoS requirements," in *Proc. ITC-15*, Washington, DC, 1997, pp. 1421–1432.
- [18] D. Mitra and I. Ziedins, "Virtual partitioning by dynamic priorities: Fair and efficient resource-sharing by several services," in *Broadband Communications: Proc. 1996 Int. Zurich Seminar Digital Commun.*, B. Plattner, Ed. New York: Springer, 1996, pp. 173–185.
- [19] M. I. Reiman and A. Schwartz, "Call admission: A new approach to quality of service," submitted for publication.
- [20] M. I. Reiman, J. Wang, and D. Mitra, "Dynamic call admission control of an ATM multiplexer with on/off sources," in *Proc. 34th IEEE Conf. Decision Contr.*, 1995, pp. 1382–1388.
- [21] J. W. Roberts, Ed., "Performance evaluation and design of multiservice networks," Final Report of the Cost 224 Project, Commission of the European Communities, Brussels, 1992.
- [22] S. M. Ross, *Applied Probability Models with Optimization Applications*. San Francisco, CA: Holden-Day, 1970.
- [23] H. Saito, "Call admission control in an ATM network using upper bound of cell loss probability," *IEEE Trans. Commun.*, vol. 40, pp. 1512–1521, 1992.
- [24] A. Schwartz and A. Weiss, *Large Deviation for Performance Analysis*. London, U.K.: Chapman & Hall, 1995.

Debasis Mitra (S'94–M'95) for a photograph and biography, see this issue, p. 678.



Martin I. Reiman received the A.B. degree in physics and math from Cornell University, Ithaca, NY, in 1974, and the Ph.D. degree in operations research from Stanford University, Stanford, CA, in 1977.

Since 1977 he has been with Bell Labs. From 1977 to 1980 he was in the Data Communications Laboratory, Holmdel, NJ. Since 1980 he has been in the Mathematical Sciences Research Center, Murray Hill, NJ. His current research interests are in the analysis, optimization, and control of stochastic service systems.

Dr. Reiman is on the Editorial Board of *Annals of Applied Probability* and *Mathematics of Operations Research*.



Jie Wang received the B.S. degree in mathematics from Beijing University, China, in 1982, the M.S. degree in applied math from the University of Alberta, Edmonton, Alta., Canada, in 1988, and the Ph.D. degree in systems engineering from the University of Pennsylvania, Philadelphia, in 1993.

He is currently Senior Member of Technical Staff at the Department of Teletraffic and Performance Analysis, AT&T Labs, Holmdel, NJ.