

# Optimal Pricing for Multiple Services in Telecommunications Networks Offering Quality of Service Guarantees

Neil J. Keon *Member, IEEE*, G. Anandalingam, *Senior Member, IEEE*

**Abstract**— We consider pricing for multiple services offered over a single telecommunications network. Each service has quality of service (QoS) requirements that are guaranteed to users. Service classes may be defined by the type of service, such as voice, video or data, as well as the origin and destination of the connection provided to the user. We formulate the optimal pricing problem as a nonlinear integer expected revenue optimization problem. We simultaneously solve for prices and the resource allocations necessary to provide connections with guaranteed QoS. We derive optimality conditions and a solution method for this class of problems, and apply to a realistic model of a multi-service communications network.

**Index Terms**— Economics, Network Design, Pricing, Quality of Service (QoS)

## I. INTRODUCTION

In this paper, we derive a nonlinear mathematical programming model for determining optimal prices for multiple services with guaranteed quality of service. We solve this model using an auction algorithm. The traditional distinction between voice networks, data networks and cable TV networks is fast becoming obsolete. In the future, multiple communications services will be available to users over a single network. These different services will vary in terms of bandwidth requirements and tolerate different quality limitations such as loss of data and delay in transmission. Our model presents a way for network providers to set prices for these services, and allocate resources such that these Quality of Service (QoS) requirements are guaranteed while expected revenues are maximized.

### A. Overview

We consider the problem of using pricing and resource allocation to manage multiple services networks with Quality-of-Service guarantees. Although the arrival rate of connection

requests can be fixed by setting the price, the observed arrivals, and thus observed revenue, are assumed to be a stochastic process. Also, the model we propose in this paper separates demand into *classes*, defined by service type as well as origin and destination. Our model maximizes *expected* revenue in order to obtain optimal prices and the resources needed to guarantee QoS for each service class.

The “resources” we consider are *buffer* space for temporarily storing users’ data traffic, and *bandwidth* for transmitting the data through the network. All connections in a particular service class (eg. voice, data, video etc.) share the same buffer and bandwidth allocation. The resource allocations must be sufficient to satisfy two QoS guarantees which are defined for every service class individually. The first QoS guarantee is a restriction on the probability of data being lost in the network. Data is lost when resources allocated to a service class based on *expected* traffic calculations are insufficient to either transmit (using bandwidth) or store (using buffer space) all the data that is actually sent to the network by all users in a service class. The second QoS guarantee is a limit on the delay, or time the data is buffered, during transmission by the network. The two QoS criteria constrain the problem in terms of buffer space and bandwidth that must be allocated to each service class for a given level of demand.

The trade-off between resource allocations among various service classes is complicated due to the fact that as the number of connections in a class increases we observe economies of scale in bandwidth allocations. This is due to two phenomena: smoothing of data traffic in buffers and statistical multiplexing gains. Economies-of-scale appear by decreasing the marginal resource allocations to any service class, i.e. by decreasing the resources required to supply one additional connection.

We construct a distributed search method called the “auction algorithm” to optimize the revenue by choosing prices, subject to constraints on the resource allocations implied by finite resources and QoS guarantees. Our search method takes advantage of several optimality properties we show for the resource allocations to each class. Using these properties we reduce the search space for the overall problem and calculate optimal solutions.

Manuscript received February 18, 2001. This research was partially funded by a grant from the National Science Foundation NCR-9612781.

Neil J. Keon is with the Cox School of Business, Southern Methodist University, Dallas TX 75275 USA (phone: 214-768-3096; fax: 214-768-4099; e-mail: nkeon@mail.cox.smu.edu).

G. Anandalingam is with R.H. Smith School of Business and the Institute for Systems Research, University of Maryland, College Park MD 20742 USA. (e-mail: ganand@rhsmith.umd.edu).

## B. Related Literature

We will briefly classify the related literature into two broad categories. First, we will discuss the bulk of the literature, which deals with pricing in best effort networks, such as the Internet. There has also been some recent work related to pricing connections for networks with QoS guarantees.

### 1) Pricing for Best Effort Service

The most celebrated packet-switched network, the Internet, offers *best effort service*, which is prone to unpredictable congestion and delays by definition. The current flow control scheme in the internet is called transmission control protocol (TCP) [2] [34]. Modifications to this scheme have been suggested that would allow multiple service classes, and induce users to behave fairly and efficiently, through simple or randomized packet marking mechanisms in the network [10][13] [14] [20]. Usage based schemes, which charge based on the actual resources used have also been proposed for these types of networks [8] [29].

Priority pricing is another suggestion for allowing multiple services over best effort networks. The most well known work discusses a second bid auction, whereby it is incentive compatible, or in the users best interests to truthfully reveal their true valuation of service in terms of a priority [22][23][24]. Another bidding paradigm, which results in a Nash equilibrium among users was suggested for the Internet in [5]. A similar approach, again requiring users to bid for service, was proposed for available bit rate service, the best-effort service offered in asynchronous transfer mode (ATM) networks [6].

A recent game theoretic model requiring users to choose from among several routes, each with its own delay, has been shown to have a stable equilibrium solution where the relative prices induce the desired operating point of the network operator [16][17]. An interesting problem of regulating the arrival of jobs presented to the network is discussed in [26], and the informational requirements are relaxed, using an adaptive on-line method in [25]. Finally, recent results have been offered, based on dynamic programming that suggest a static price schedule results in network performance that is near optimal compared to congestion dependent pricing [27].

### 2) Pricing with QoS Guarantees

Typically a network with guaranteed QoS, such as the current voice network, must employ a call admission policy in order to satisfy the guarantees to users. In an alternative approach to call admission, it has been suggested that users guarantee their own QoS by purchasing the required bandwidth and buffer resources for their desired QoS directly from the network [19] [21]. In [31], an analysis of a market-based methodology is offered as evidence that pricing schemes can offer efficient resource allocations in connection-oriented networks offering QoS guarantees. (See also [33]) These models feature a conventional view of pricing per connection and announcing the price to the public. The paper by de Veciana et al. [32] combines the conventional practice of prices being fixed and announced, with the alternative view of allowing users to guarantee their own QoS by purchasing

bandwidth and buffers directly. Resources are shared among users, incorporating multiplexing gains. On-line negotiation is also suggested as a framework for connection set-up and allocation of network resources, using effective bandwidth as a base for pricing in [11] [12]. In a related paper that deals with single link point-to-point ATM networks, the pricing problem was formulated as a constrained optimal control problem, and solved using a three-stage solution procedure [35]. The limitation of the approach is computational tractability.

Our paper belongs to this class of work, but uses a novel optimization model based on nonlinear mathematical programming. We use an optimal pricing scheme to moderate the demand for connection requests from different classes of service. Each "class" of service is defined by both the type of service (eg. Voice, data, video), and by where the telecommunications flow originates, and where it is destined to. Thus, our paper tries to take into account the use of resources in the entire network. Based on the connection requests for each class of service at the specified prices, the model estimates optimal buffer and bandwidth allocations to satisfy QoS requirements. Another contribution of this paper is the algorithm we propose to obtain the optimal solutions.

## C. Organization of the Paper

This paper is organized as follows: In Section II, we will describe the Network Model that we are considering, and examine three important features, call blocking, loss probability, and delay in some detail. Section III contains the complete formulation of the expected revenue maximization model, and presents optimality properties. Section IV provides details of the Auction Algorithm used to solve the problem. We present the implementation of the algorithm to an extensive example in Section V. We end the paper with concluding remarks in Section VI.

## II. A MODEL FOR OPTIMAL PRICING WITH QoS GUARANTEES

In this section, we will give an overview of the model that we will be using to derive optimal prices for the multiple services in a telecommunications network where quality of service (QoS) is guaranteed. We will first describe an example network that will help us focus our thoughts on the particular features that we model, present the characteristics of the model, and then derive the particular mathematical constructs of these characteristics. In section 3, we will present the complete mathematical model, and in Section 4, we will present details of the algorithm to compute the optimal multi-service prices, based on the optimality conditions of the mathematical model.

### A. The Network

We will use the ring type network shown in Fig. 1, only to *highlight* the particular features that we are modeling. Note, however, that all features we consider in our model are found in any type of network. At each switch there are inputs and outputs of telecommunications traffic. Traffic inputs are from

many different *types* of service that either originate at that switch, or else arrive at the switch from other originating nodes through an incoming bandwidth pipe. Traffic outputs from the switch either go along an outgoing bandwidth pipe, or directly to a user connected to that switch.

Without loss of generality with respect to the network characteristics, we assume that buffer space is used at originating switches to collect and smooth data traffic, subject to delay constraints. The data traffic is then transported in bandwidth pipes through the rest of the path to its destination. We allocate aggregate quantities of buffer and bandwidth to connections grouped by service type,  $i$ , origin,  $j$ , and destination,  $k$ . A service *class* is therefore denoted by the triple,  $(i,j,k)$ . Each service type  $i$  has traffic characteristics defined by *average* bandwidth  $r_i$  and *peak* bandwidth  $R_i$ . For example, voice and video have different average and peak bandwidth requirements.

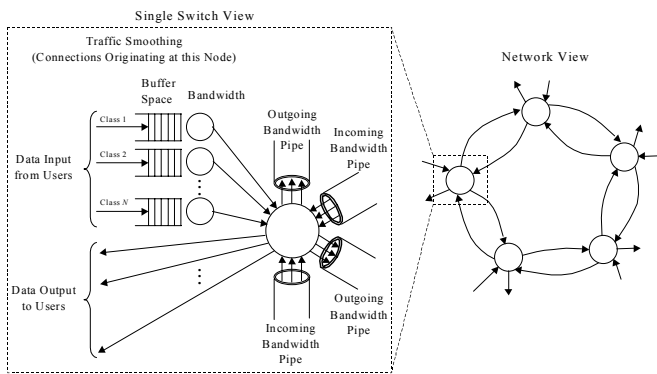


Fig. 1. Network (Ring) Service Model at an Individual Switch.

### B. Overview of the Model

The model is derived from the point-of-view of the service provider. The service provider sets prices for the different multiple services in order to maximize his/her revenue subject to a set of constraints and conditions that are necessary to ensure both quality of service and flow balance. The problem that the service provider solves can be summarized as:

- *Maximize Expected Revenue*: This is a product of prices charged for each service and the demand for each service, taken to be an arrival rate of connection requests, *that can be accommodated* at that price.
- Subject to the following constraints:
- *Limited Capacity*: The network switches have limited capacity (bandwidth)
- *Limit the Blocking of Connection Requests*: Given bandwidth capacity constraints, there is a limited number of connections that can be supported in order to guarantee QoS, and some requests may have to be blocked. We set a limit on the blocking probability for each service class. Note that the capacity constraint and the limitation on blocking probability will in turn affect the number of connection requests that can be accommodated.
- *Limit the Probability of Packet Loss*: Different services can tolerate different levels of packet loss and still

guarantee QoS. We set a limit on the “equivalent capacity” (defined in detail later) allocated to each service in order to limit the probability of packet loss.

- *Limit the Maximum Delay*: Different services can tolerate different levels of delay and still guarantee QoS. We set a limit on the maximum allowable delay for each service.

We will now describe each component of the model in greater detail using quantities illustrated in Fig. 1. We will present the complete mathematical formulation of the problem in section II, and derive optimality conditions.

### C. Deriving Demand

The quantity demanded by users,  $\lambda_{ijk}$ , is taken to be the arrival rate of connection requests for service  $i$ , with the connection originating at switch  $j$  and terminating at switch  $k$ . This quantity is determined by a demand function  $f_{ijk}(p_{ijk})$  where  $p_{ijk}$  is the price charged per unit of time the connection is open. Based on current empirical results [1][18], in our model we use the demand function below:

$$\lambda_{ijk} = a_{ijk} p_{ijk}^{-\epsilon_{ijk}} \quad (1)$$

where  $\epsilon_{ijk}$  is the (constant) elasticity of demand for service class  $(i, j, k)$ . One important property of (1) is that marginal revenue with respect to prices will always be negative (meaning lower prices always increase revenue) provided that  $\epsilon_{ijk} > 1$ . We assume elasticities in excess of unity for the remainder of the paper. Again, this is suggested by available research [1][18].

### D. Call-Blocking and Expected Revenue

We are modeling a network offering connections with guaranteed QoS. Under these conditions, there is a limited number of connections that can be supported, and some requests may have to be blocked if resources are already reserved for other requests. We would like to minimize the blocking of connection requests, and capture this in our optimization model as a constraint.

For each service class, denoted by  $(i,j,k)$ , we model the number of ongoing connections as an  $M/GI/m/m$  queue. The capacity in the queuing model,  $m$ , refers to the maximum number of connections that we will allow to be admitted to a given service class, a variable we define as  $NCMAX_{ijk}$ . For consistency with this model, we state a number of assumptions: The arrival rate of connection requests, for a service class, given by  $\lambda_{ijk}$  (defined above) characterizes inter-arrival times that are independently and identically distributed exponential random variables. The average holding time of a connection,  $T_{ijk}$ , is known, and the holding times of individual connections can occur according to any general distribution, but are independently and identically distributed for each connection. The traffic intensity of a class,  $\rho_{ijk}$ , is the product of arrival rate and average holding time, i.e.,  $\rho_{ijk} = \lambda_{ijk} T_{ijk}$ .

The properties of the  $M/GI/m/m$  queuing model are well known [4]. The probability of blocking a request within any service class,  $BP_{ijk}$ , is given by the Erlang B formula:

$$\begin{aligned}
BP_{ijk} &= P(NC_{ijk} = NCMAX_{ijk}) \\
&= \frac{\rho_{ijk}^{NCMAX_{ijk}}}{NCMAX_{ijk}!} \bigg/ \sum_{n=0}^{NCMAX_{ijk}} \frac{\rho_{ijk}^n}{n!}
\end{aligned} \quad (2)$$

If the number of open connections,  $NC_{ijk}$ , is equal to the maximum number permitted for each class,  $NCMAX_{ijk}$ , when a user's request for service is received, the connection will be blocked and lost to the system. Otherwise (i.e.  $NC_{ijk} < NCMAX_{ijk}$ ) the request will be admitted.

The expected number of busy connections is given by:

$$E[NC_{ijk}] = \rho_{ijk} (1 - BP_{ijk}) \quad (3)$$

We define the expected rate of revenue generation attributable to a service class,  $REV_{ijk}$ , as the price charged per unit time for a single connection,  $p_{ijk}$ , multiplied by the expected number of ongoing connections,  $E[NC_{ijk}]$ :

$$REV_{ijk} = p_{ijk} E[NC_{ijk}] \quad (4)$$

### E. Probability of Loss and Equivalent Capacity

In the previous section, we modeled the connection level traffic so we could quantify blocking probability for our model. Now we model the packet level traffic for the same set of connections, so we can quantify one of the QoS parameters, namely probability of lost packets. This will also allow us to relate the service classes to the necessary resource assignments, buffer space and bandwidth, to satisfy the QoS. At the switch, packets from all connections of the same class ( $i,j,k$ ) share a FIFO queue of size  $B_{ijk}$ . Up to  $NCMAX_{ijk}$  connections, with statistically identical data traffic, may share the allocated buffer at any time. Packet loss,  $PLOSS_{ijk}$ , occurs when there is buffer overflow.

To ensure that the guaranteed loss probability is satisfied for each service class we require that the minimum bandwidth allocation per offered connection be at least equal to the *equivalent capacity*. Equivalent capacity,  $c$ , for a single source (connection) is defined as, "the service rate of the queue that corresponds to a given (packet) loss,"  $PLOSS_{ijk}$  [9][28]. If bandwidth in excess of the equivalent capacity is assigned to a connection, the observed packet loss is less than  $PLOSS_{ijk}$ . There are two main effects that determine the equivalent capacity of a single source. *Effective bandwidth* refers to the fact that smoothing of data traffic in the buffer reduces the bandwidth required to achieve a prescribed loss rate of packets. *Multiplexing gains* refers to the gain in bandwidth required due to the statistical effect of mixing the data traffic of independent connections. We use the term *equivalent capacity* when referring to the combination of these two effects on bandwidth assignment (per connection) to a number of connections with guaranteed QoS. We now state some general properties of equivalent capacity.

In our model, all connections of a service class ( $i,j,k$ ) share a single FIFO buffer. Consider a connection of a given service type,  $i$ , which is distinct from service "class" where we include origin-destination information as well. The user is either sending data at a *peak* transmission rate,  $R_i$ , (i.e. the "on" state) or is sending no traffic at all (i.e. the "off" state).

The *average* rate of data transmission for the individual connection is given by  $r_i$ . The on and off periods are assumed to be exponentially distributed, with the average length of an on periods given by  $b_i$ . QoS requirements for this type of data traffic model have been well studied, e.g. [9]. The traffic from all connections of the same type,  $i$ , is statistically identical, i.e., all connections classified as the same type have the same parameters  $r_i$ ,  $R_i$ , and  $b_i$ .

Clearly, the equivalent capacity must be greater than or equal to the mean rate of data traffic,  $r_i$ , and less than or equal to the peak rate of data traffic,  $R_i$ .

$$r_i \leq c(NCMAX_{ijk}, B_{ijk}, PLOSS_{ijk}) \leq R_i \quad (5)$$

As figure 2 shows equivalent capacity is decreasing in the amount of allocated buffer space,  $B_{ijk}$ .

$$\frac{\partial c(NCMAX_{ijk}, B_{ijk}, PLOSS_{ijk})}{\partial B_{ijk}} \leq 0 \quad (6)$$

As the number of connections for which we allocate resources,  $NCMAX_{ijk}$ , increases, equivalent capacity will decrease reflecting multiplexing gains.

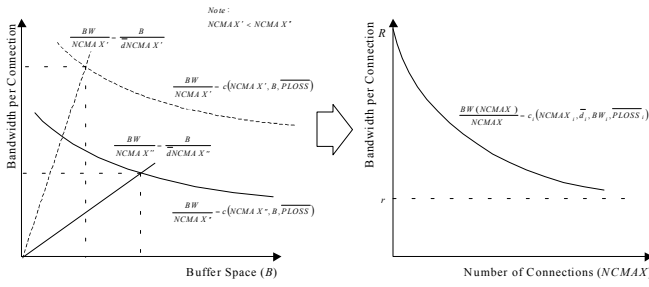
$$\frac{\partial c(NCMAX_{ijk}, B_{ijk}, PLOSS_{ijk})}{\partial NCMAX_{ijk}} \leq 0 \quad (7)$$

Finally, as we allow greater loss probabilities,  $PLOSS_{ijk}$ , the equivalent capacity will decrease.

$$\frac{\partial c(NCMAX_{ijk}, B_{ijk}, PLOSS_{ijk})}{\partial PLOSS_{ijk}} \leq 0 \quad (8)$$

An in-depth analysis of equivalent capacity is outside the focus of this paper. Detailed discussions of buffering and multiplexing gains are given in [2][9][28][32] and [34]. However, it should be noted that our approach remains unchanged for *any* derivation of equivalent capacity, which satisfies the general properties for equivalent capacity, given in (5) – (8) above. The inequalities admit special cases such as constant bit rate traffic, with no smoothing or multiplexing gains. As we shall see in the next section, the data traffic model and equivalent capacity expressions in (5) – (8) are sufficient for formulating packet loss in our revenue-maximizing model.

It will become apparent later that we wish to solve buffer assignment and bandwidth assignment simultaneously. Therefore, we now illustrate some extensions of the equivalent capacity properties when maximum buffer delay is held constant, i.e. the buffer is sized according to a constant maximum delay given by  $d = B/BW$ . Fig. 2 shows the general relationship between equivalent capacity (average bandwidth per connection) and number of connections when maximum delay in the buffer is held fixed. This relationship is derived based on (5) – (8) above (the service class indices, ( $i,j,k$ ) are omitted for clarity.) Note that equivalent capacity, as illustrated in Fig. 2 below, gives the bandwidth required to meet two QoS criteria, namely loss probability,  $PLOSS_{ijk}$ , and maximum delay,  $d_{ijk}$ .



**Fig. 2. Illustration of Average Bandwidth Allocated per Connection, Using Mapping from NCMAX to BW.**

The average bandwidth allocated per connection, as a result of the mapping from  $B$  to  $BW$  given by (5) – (8) is illustrated by the downward sloping curves in the left panel above. An increase in the number of connections served, from  $NCMAX'$  to  $NCMAX''$  results in the downward shift the the average bandwidth assignment per connection. The straight lines represent the relationship bewteen buffer,  $B$ , and bandwidth  $BW$ , for fixed maximum delay, i.e.  $BW = dxB$ . Since the figure is drawn in a space representing average bandwidth per connection,  $BW/NCMAX$ , the shift outward in the slope of the lines results from the same increase in the number of connections served, from  $NCMAX'$  to  $NCMAX''$ . Collecting the solutions to such a system of two equations in two unknowns yields the downward sloping curve in the right panel of Fig. 2.

The figure above is conceptual. In our numerical example, we will use the equivalent capacity results from [9] to calculate the equivalent capacity per connection summarized below:

$$c(NCMAX_{ijk}, B_{ijk}, PLOSS_{ijk}) = \min\left(\hat{c}, \frac{\mu + \alpha'\sigma}{NCMAX_{ijk}}\right) \quad (9)$$

where,

$$\hat{c} = \frac{\alpha R_i - B_{ijk} + \sqrt{(\alpha R_i - B_{ijk})^2 + 4\alpha B_{ijk} r_i}}{2\alpha} \quad (10)$$

$$\alpha = \ln(1/PLOSS_{ijk}) b_i \left(1 - \frac{r_i}{R_i}\right) \quad (11)$$

$$\mu = r_i NCMAX_{ijk} \quad (12)$$

$$\sigma = \sqrt{r_i (R_i - r_i) NCMAX_{ijk}} \quad (13)$$

$$\alpha' = \sqrt{-2 \ln(PLOSS_{ijk}) - \ln(2\pi)} \quad (14)$$

The equivalent capacity per connection, (9), is calculated as the minimum of two distinct approximations. The effective bandwidth is approximated by (10). The second term in the minimum expression in (9) reflects multiplexing gains. This approximation is based on the stationary bit rate. To calculate this expression, we need the additional expression (11), as well as the mean of the aggregate bit rate, (12), the standard deviation of the aggregate bit rate, (13), and an approximate inversion of the normal distribution, (14). The calculations above separate equivalent capacity into regions dominated by either smoothing effects in the buffer or multiplexing gains.

These expressions are one example of suggested computational methods for equivalent capacity that satisfy (5) – (8).

#### F. Delay

The second QoS parameter that we consider is the maximum allowable delay,  $d_{ijk}$  for each service class. We ignore transmission delay and consider only delay in the buffer. Therefore we can set bounds for the potential delay a packet may experience quite easily:

$$d_{ijk} \leq \frac{B_{ijk}}{BW_{ijk}} \quad (15)$$

where  $B_{ijk}$  is the buffer space and  $BW_{ijk}$  is the bandwidth allocated to service (i, j, k). The maximum delay any packet may experience, given by (15), is simply the size of the allocated buffer space divided by the allocated bandwidth. The buffer is served on a first-in first-out (FIFO) basis.

### III. THE OPTIMIZATION MODEL

#### A. The Model

Given the derivations above, we are now in a position to present the complete mathematical formulation of the revenue maximization model that would yield optimal prices and resource allocations. We wish to solve for the price,  $p_{ijk}$ , for each service class denoted by (i,j,k), as well as the volume of service offered,  $NCMAX_{ijk}$ . The volume of service offered refers to how many connections we can serve simultaneously, given the bandwidth,  $BW_{ijk}$ , reserved all along the path between j and k, and buffer space,  $B_{ijk}$ , reserved at the origin switch. We assume that bandwidth is limited but buffer space is not. The network service is assumed to be offered using a connection admission policy, so that the QoS requirements for probability of loss,  $PLOSS_{ijk}$ , and delay,  $d_{ijk}$ , are satisfied within certain limits. We include the blocking probability,  $BP_{ijk}$ , resulting from the use of a connection admission policy to ensure that we have sufficient resources set aside for service class (i,j,k), such that a satisfactory proportion of connection requests can be admitted when they are requested.

The total expected revenue from all service connections in the network is given by  $\sum_{i=1}^N \sum_{j=1}^{NS} \sum_{k=1}^{NS} (p_{ijk} E\{NC_{ijk}\})$ . Because the expected number of connections arrivals  $NC_{ijk}$  is a function of the arrival rate, the blocking probability, and the average holding time of a connection, we can replace the objective function with  $\sum_{i=1}^N \sum_{j=1}^{NS} \sum_{k=1}^{NS} (p_{ijk} f\{\lambda_{ijk}, BP_{ijk}, T_{ijk}\})$ .

Thus the revenue optimization model that will yield optimal prices and resource allocation is given by (PNet)

$$Max_{p_{ijk}, B_{ijk}, BW_{ijk}, NCMAX_{ijk}} \sum_{i=1}^N \sum_{j=1}^{NS} \sum_{k=1}^{NS} (p_{ijk} f\{\lambda_{ijk}, BP_{ijk}, T_{ijk}\}) \quad (16)$$

subject to,

$$BP_{ijk}(\lambda_{ijk}(p_{ijk}), T_{ijk}, NCMAX_{ijk}) \leq \overline{BP}_{ijk}, \quad (17)$$

$$1 \leq i \leq N, 1 \leq j \leq NS, 1 \leq k \leq NS$$

$$c(NCMAX_{ijk}, B_{ijk}, \overline{PLOSS}_{ijk}) \leq \frac{BW_{ijk}}{NCMAX_{ijk}}, \quad (18)$$

$$1 \leq i \leq N, 1 \leq j \leq NS, 1 \leq k \leq NS$$

$$\frac{B_{ijk}}{BW_{ijk}} \leq \overline{d}_{ijk}, \quad 1 \leq i \leq N, 1 \leq j \leq NS, 1 \leq k \leq NS \quad (19)$$

$$\sum_{i=1}^N \sum_{j=1}^{NS} \sum_{k=1}^{NS} BW_{ijk} IVP_{ijkx} \leq \overline{BW}_x, \quad 1 \leq x \leq NS \quad (20)$$

$$IVP_{ijkx} = \begin{cases} 1 & , \text{if } x \in VP_{ijk} \\ 0 & , \text{otherwise} \end{cases}, \quad (21)$$

$$1 \leq i \leq N, 1 \leq j \leq NS, 1 \leq k \leq NS, 1 \leq x \leq NS$$

$$\lambda_{ijk} \geq 0, p_{ijk} \geq 0, BW_{ijk} \geq 0, B_{ijk} \geq 0 \quad (22)$$

$$NCMAX_{ijk} \geq 0, \text{ integer} \quad (23)$$

where,

$REV_{ijk}$  = rate of revenue generation associated with service class  $(i,j,k)$

$c(\cdot)$  = equivalent capacity of a single connection

$VP_{ijk}$  = {set of all  $x$  in the path for class  $i$  originating at  $j$  and terminating at  $k$ ,  $1 \leq x \leq NS$ }

$\overline{BP}_{ijk}$  = maximum blocking probability for a connection from service class  $(i,j,k)$

$\overline{PLOSS}_{ijk}$  = maximum packet loss probability for a for a connection from service class  $(i,j,k)$

$\overline{d}_{ijk}$  = maximum delay for a for a connection from service class  $(i,j,k)$

$\overline{BW}_j$  = the capacity (bandwidth) at switch  $j$

In this formulation, the objective function, (16), seeks to maximize the average rate of revenue generation from ongoing connections, which is given by price multiplied by the *expected* number of connections. The constraints restrict the performance measures and resource assignments to be within bounds set outside the problem as *network policy*. Budgets on each quantity are denoted by a bar overhead. Constraint (17) restricts the call-blocking probability for every class below some prescribed limit. Constraint (18) ensures the probability of loss for each service class,  $\overline{PLOSS}_{ijk}$ , is satisfied. The buffer delay for each service class is constrained to be less than or equal to the limit given by the QoS guarantee in constraint (19). We have a bandwidth capacity constraint at each switch, (20), but we assume there is no capacity on allocated buffer space, and there are no link capacities. The capacity at each switch must accommodate all traffic originating at the switch as well as traffic originating elsewhere but routed through the switch. We assume the routes are known and fixed. There is an indicator function, (21), for every class  $(i,j,k)$ , which indicates if any switch is included in the path. There are a number of non-negativity

constraints, given by (22). Finally there is an integrality constraint on the maximum number admitted to each class,  $NCMAX_{ijk}$ , given by (23), since connections can only be admitted in discrete quantities.

We wish to call attention to one limitation of our formulation of the capacity constraints, which may over-assign bandwidth along the paths followed. The effective bandwidth of the aggregate traffic for the service class  $(i,j,k)$  may be reduced after smoothing in the buffer at the originating switch. Our formulation assigns bandwidth all along the path based on the effective bandwidth at the originating switch and does not take into account the statistical properties of the output traffic at the originating switch. The simplicity of the formulation outweighs our concern regarding excess bandwidth assignments.

### B. Optimality Properties

In general, the problem given in (PNet) is a nonlinear, non-convex, mixed integer problem. However, there are a number of necessary conditions for an optimal solution, which we will use to search for a solution. First we will discuss necessary conditions related to constraints (17) – (19) in problem (PNet). These constraints are relatively simple as they apply to each service class independently of all other service classes. Based on this first set of necessary optimal conditions, we will then state a pair of optimality conditions reflecting the marginal values of each service class  $(i,j,k)$  in the solution.

#### 1) Resource Allocation to Individual Service Classes

There are a number of conclusions we can immediately draw about optimal solutions to the problem (PNet). First we consider the resources that must be assigned to each individual service class.

#### **Theorem:**

Given downward-sloping demand curves and plentiful buffer space at each originating switch, if an optimal solution to (PNet) does not coincide with the condition that marginal revenue equal to zero, i.e. the partial derivatives of (1) with respect to prices are all less than zero at optimality, then the optimal solution to (PNet) must satisfy the following properties:

$$BP_{ijk}(\lambda_{ijk}(p_{ijk}), T_{jk}, NCMAX_{ijk}) = \overline{BP}_{ijk}, \quad (24)$$

$$1 \leq i \leq N, 1 \leq j \leq NS, 1 \leq k \leq NS$$

$$\frac{BW_{ijk}}{NCMAX_{ijk}} = c_{ijk}(NCMAX_{ijk}, B_{ijk}, \overline{PLOSS}_{ijk}), \quad (25)$$

$$1 \leq i \leq N, 1 \leq j \leq NS, 1 \leq k \leq NS$$

$$\left( (NCMAX_{xjk} + 1) c_{xjk}(NCMAX_{xjk} + 1, B_{xjk}, \overline{PLOSS}_{xjk}) + \sum_{\substack{i=1 \\ i \neq x}}^N \sum_{k=1}^{NS} NCMAX_{ijk} c_{ijk}(NCMAX_{ijk}, B_{ijk}, \overline{PLOSS}_{ijk}) \right) > \overline{BW}_j$$

$$1 \leq x \leq N, 1 \leq j \leq NS \quad (26)$$

$$B_{ijk} = \overline{d}_{ijk} BW_{ijk}, \quad 1 \leq i \leq N, 1 \leq j \leq NS, 1 \leq k \leq NS \quad (27)$$

**Proof:** See Appendix for proof.

The call-blocking probability for all service classes is set to its greatest permitted value for all service classes in (24). Incidentally, this is the only condition, which pertains directly to prices. The bandwidth assigned to each service class, for the maximum number of connections admitted, will be equal to the corresponding equivalent capacity, as stated in (25). The third condition, (26), states that because connections have to be integer, allocation of the equivalent capacity for any additional connections would violate feasibility. This means that no further increases in bandwidth assignments are possible at an optimal (feasible) solution. Buffer space is assigned for all services originating at each switch, such that packets for each service may experience a delay up to the tolerated delay, (27), assuming that buffer space at each switch is plentiful.

The optimal resource allocation is to assign the appropriate effective bandwidth to each service class, (25), and buffer space proportional to the equivalent capacity (27). Typically, buffer space is thought of as a parameter in the equivalent capacity calculation, while we are treating it as a variable (along with bandwidth), for which we solve a system of two equations in two unknowns. That is, we solve (25) and (27) for  $BW_{ijk}$  and  $B_{ijk}$ :

$$BW_{ijk}(NCMAX_{ijk}) = BW : \frac{BW}{NCMAX_{ijk}} = c(NCMAX_{ijk}, \bar{d}_{ijk} BW_{ijk}, \overline{PLOSS}_{ijk}) \quad (28)$$

Going back to Fig. 2 which shows equivalent capacity as a decreasing function of buffer space, and the number of connections served, the optimality property, (27), is illustrated with a linearly increasing function of buffer allocation, where both sides of the expression have been divided by  $NCMAX_{ijk}$ . As the number of connections served increases, the slope of the line becomes less steep. Based on this simple graphical analysis, the system of equations, (25) and (27), must have a unique solution and be decreasing in terms of  $NCMAX_{ijk}$ . Note that if there are no smoothing effects of multiplexing gains, as with constant bit rate traffic ( $R_i = r_i$ ), then the allocation per connection is simply a constant value. In all other cases, the total allocation must therefore be increasing but the allocation per connection may be decreasing, depending on the properties of the service class. This shows economies of scale in the optimal resource allocation.

Using another optimality property from Theorem 1, the call-blocking constraint will be binding by (24). We can solve for the the optimal arrival rate,  $\lambda_{ijk}$ , for a given  $NCMAX_{ijk}$ , using (1):

$$\rho_{ijk}^*(NCMAX_{ijk}) = \rho : \overline{BP}_{ijk} = \frac{\rho^{NCMAX_{ijk}}}{NCMAX_{ijk}!} \bigg/ \sum_{n=0}^{NCMAX_{ijk}} \frac{\rho^n}{n!} \quad (29)$$

$$\lambda_{ijk}(NCMAX_{ijk}) = \frac{\rho_{ijk}^*}{T_{ijk}} \quad (30)$$

The arrival rate,  $\lambda_{ijk}$ , is determined by the limit on call-blocking. We must first calculate the maximum traffic intensity for  $NCMAX_{ijk}$ , given in (29), and then calculate the

arrival rate, using (30). Because the resource allocation tables are independent of demand we have not yet related the arrival rate to a price. The bid tables, presented in the next section, will relate the price required,  $p_{ijk}$ , to produce the desired arrival rate,  $\lambda_{ijk}$ , as well as the marginal valuation at the volume of service provided,  $NCMAX_{ijk}$ .

## 2) Resource Allocation Trade-Offs Between Service Classes

The optimality properties, from sub-section 1) above, reduce the problem to choosing the optimal values for  $NCMAX_{ijk}$ . For the time being, consider a *linear relaxation* to (PNet), eliminating the integrality constraints on  $NCMAX_{ijk}$ . For solutions satisfying the optimality properties given by (24) – (27), the Karush-Kuhn-Tucker necessary conditions are trivially satisfied, with the following exception, which results from differentiating the set of constraints (20):

$$\frac{\partial REV_{ijk}}{\partial NCMAX_{ijk}} = \sum_{x=1}^{NS} v_x \frac{\partial BW_{ijk}}{\partial NCMAX_{ijk}} IVP_{ijkx}, \quad (31)$$

$$1 \leq i \leq N, 1 \leq j \leq NS, 1 \leq k \leq NS$$

Recall that  $IVP_{ijkx}$  in (31) was defined as an indicator parameter in the problem, (PNet). For each service class  $(i,j,k)$ , we divide both sides of (31) by  $\partial BW_{ijk} / \partial NCMAX_{ijk}$ , and simplify the necessary optimality conditions:

$$u_{ijk} = \sum_{x=1}^{NS} v_x IVP_{ijkx}, \quad 1 \leq i \leq N, 1 \leq j \leq NS, 1 \leq k \leq NS \quad (32)$$

Economic interpretation of (32) is as follows: Any service class must yield a marginal return per unit of bandwidth equal to the sum of the marginal values for all switches ( $v_x$  for switch  $x$ ) along the route for the given class  $(i,j,k)$ .

Incorporating the integrality requirement in  $NCMAX_{ijk}$ , we define marginal values per unit of bandwidth for increasing or reducing the number of connections in the solution,  $u_{ijk}^+$  and  $u_{ijk}^-$  respectively:

$$u_{ijk}^+(NCMAX_{ijk}) = \frac{REV_{ijk}(NCMAX_{ijk} + 1) - REV_{ijk}(NCMAX_{ijk})}{BW_{ijk}(NCMAX_{ijk} + 1) - BW_{ijk}(NCMAX_{ijk})} \quad (33)$$

$$u_{ijk}^-(NCMAX_{ijk}) = \frac{REV_{ijk}(NCMAX_{ijk}) - REV_{ijk}(NCMAX_{ijk} - 1)}{BW_{ijk}(NCMAX_{ijk}) - BW_{ijk}(NCMAX_{ijk} - 1)} \quad (34)$$

The marginal valuations,  $u_{ijk}^-$  from (33) or  $u_{ijk}^+$  from (34), are simply the changes in expected revenue divided by the change in bandwidth allocation, for an increase or decrease of one in  $NCMAX_{ijk}$ . Note that by definition,  $u_{ijk}^+(NCMAX_{ijk})$  equals  $u_{ijk}^-(NCMAX_{ijk} + 1)$ .

Discrete approximations of the continuous necessary optimality conditions (32) are:

$$u_{ijk}^+(NCMAX_{ijk}) \leq \sum_{x=1}^{NS} v_x IVP_{ijkx} \leq u_{ijk}^-(NCMAX_{ijk}), \quad (35)$$

$$1 \leq i \leq N, 1 \leq j \leq NS, 1 \leq k \leq NS$$

In stating (35), we assume that marginal revenue is

decreasing in  $NCMAX_{ijk}$ , i.e.  $u_{ijk}^- > u_{ijk}^+$ . By (35) it is not profitable to change the bandwidth allocated to any service class. The marginal value from an increased allocation is less than the sum of marginal values of bandwidth at switches along the path.

#### IV. NETWORK AUCTION ALGORITHM

The search procedure for the global network solution seeks to identify the marginal valuations of bandwidth at all switches,  $v_j$ , that will maximize expected revenue. Based on the optimality conditions (24) to (26), it is clear that, the search will have to test and adjust marginal values at each switch until we find a solution that allocates bandwidth fully at each switch in order to satisfy the optimality conditions.

##### A. Storage of Problem Data in Bid Tables

We can exploit the first set of necessary optimality properties above, (24) – (27) to simplify the search space for the the problem (PNet). These properties dictate that optimal allocations of  $BW_{ijk}$  can be calculated from (28) based on the values of  $NCMAX_{ijk}$ . There is a unique arrival rate of connection requests,  $\lambda_{ijk}$ , associated with  $NCMAX_{ijk}$ , according to the call-blocking property, (29). In turn prices,  $p_{ijk}$ , are related to values of  $NCMAX_{ijk}$  through the arrival rate,  $\lambda_{ijk}$ , given by the demand functions, (1). This gives us a reduced space of problem data, where all other problem variables are calculated as functions of  $NCMAX_{ijk}$ , which are integer-valued variables. Furthermore, we have defined marginal values of service classes,  $u_{ijk}^-$  and  $u_{ijk}^+$ , for any value of  $NCMAX_{ijk}$ , in (33) and (34).

We can calculate all the variables referred to above off-line and summarize the search space of the problem in tables indexed by the integer-valued variables,  $NCMAX_{ijk}$  (Table 1). We call Table I a “Bid Table” for the following reason: The marginal values  $u_{ijk}^+$  (or  $u_{ijk}^-$ ) represent the maximum amount a rational agent would bid (or accept) per unit of bandwidth to add (or remove) a connection of service class  $(i,j,k)$ . Note that the bid table contains much less information than all the feasible values of  $NCMAX_{ijk}$ ,  $BW_{ijk}$ ,  $B_{ijk}$  and  $p_{ijk}$ . thus making it much faster to solve problem (PNet).

TABLE I  
DEFINITION OF A BID TABLE

$NCMAX_{ijk}$	$BW_{ijk}$ eq. (28)	$B_{ijk}$ eq. (27)	$\lambda_{ijk}$ eq. (30)	$p_{ijk}$ eq. (1)	$u_{ijk}^+$ eq. (33)	$u_{ijk}^-$ eq. (34)
1	$BW_{ijk}(1)$	$B_{ijk}(1)$	$\lambda_{ijk}(1)$	$p_{ijk}(\lambda_{ijk}(1))$	$u_{ijk}^+(1)$	$u_{ijk}^-(1)$
2	$BW_{ijk}(2)$	$B_{ijk}(2)$	$\lambda_{ijk}(2)$	$p_{ijk}(\lambda_{ijk}(2))$	$u_{ijk}^+(2)$	$u_{ijk}^-(2)$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$

##### B. Retrieval of Problem Data from Bid Tables

For any given set of marginal bandwidth valuations,  $v_j$ , we look up the “bid” and the associated number of connections and resource allocations. Strictly speaking, we select  $NCMAX_{ijk}$  for a given  $v_j$  according to the following rule:

$$NCMAX_{ijk}(v_1, v_2, \dots, v_{NS})$$

$$= \max \left( NCMAX_{ijk} : u_{ijk}^+ < \sum_{i=1}^{NS} v_x IVP_{ijkx} < u_{ijk}^- \right) \quad (36)$$

We search from the bottom of the bid tables and take the first (and largest) value of  $NCMAX_{ijk}$  for which  $u_{ijk}^+ < v_j < u_{ijk}^-$ . For relatively small values of  $NCMAX_{ijk}$  the  $u^-/u^+$  values may be increasing due to large multiplexing gains. Since we are interested in revenue maximization, (36) selects the largest value of  $NCMAX_{ijk}$ , which satisfies the optimality property, (35). The simplest way to find this value is by looking up the bid table data starting from the bottom of the table. We will now present a method called the “auction algorithm” for ensuring that the  $NCMAX$  we choose satisfies the optimality conditions.

##### C. Bounds on Marginal Value of Bandwidth

For convenience in describing the search procedure, we define bounds on the *optimal set of marginal valuations* of bandwidth for all switches:

$$\{\underline{v}_1, \underline{v}_2, \dots, \underline{v}_{NS}\} < \{v_1^*, v_2^*, \dots, v_{NS}^*\} < \{\bar{v}_1, \bar{v}_2, \dots, \bar{v}_{NS}\} \quad (37)$$

where,

$v_j^*$  = optimal marginal valuation at switch  $j$

$\underline{v}_j$  = a lower bound on the value of  $v_j^*$  at switch  $j$

$\bar{v}_j$  = an upper bound on the value of  $v_j^*$  at switch  $j$

Clearly, using the bid table and (37), we can guarantee a *set of marginal values* is either an upper bound or a lower bound on an optimal set of marginal values:

$$\sum_{i=1}^N \sum_{j=1}^{NS} \sum_{k=1}^{NS} BW_{ijk}(\underline{v}_1, \underline{v}_2, \dots, \underline{v}_{NS}) IVP_{ijkx} > \overline{BW}_x \quad \forall x \quad (38)$$

$$\sum_{i=1}^N \sum_{j=1}^{NS} \sum_{k=1}^{NS} BW_{ijk}(\bar{v}_1, \bar{v}_2, \dots, \bar{v}_{NS}) IVP_{ijkx} < \overline{BW}_x \quad \forall x \quad (39)$$

A set of marginal values that undervalues the bandwidth at all switches results in an over-assignment of bandwidth at all switches and be infeasible. Thus, we can use such a set of marginal values as a *lower bound* on the set of optimal marginal values. Likewise a set of marginal values which overvalues the bandwidth at all switches will be feasible and can be used as an upper bound on an optimal set of marginal values. Note that while we offer bounds on the optimal value of  $v_j^*$  in (37), we claim only to be seeking bounds for a local optimal solution to the problem (PNet), which is non-convex and contains integer-valued variables.

##### D. Search Procedure

Our search method begins with arbitrary bounds on the optimal marginal valuation of bandwidth at each switch, which satisfy (38) and (39). Each iteration, we decrease the Euclidean distance between the two sets of bound and change the direction of the line segment between the two sets of bounds in the marginal value space. When the upper and lower bounds are very near to each other, and we have a feasible solution that fully assigns capacity at all switches, we terminate the search and obtain a near optimal solution, where

capacity is fully assigned and revenue is very close to optimal. The infeasible solution given by the lower bounds on marginal value of bandwidth at each switch will also provide a bound on the distance from a local optimal solution.

The search algorithm is as follows:

**Step 1 Initialization:** For initialization, we simply require an initial set of bounds on the optimal marginal values as defined in (38) and (39).

**Step 2 Identifying Under-Valued (and Over-Valued) Bandwidth:** We wish to identify the best feasible and worst infeasible solutions that lie along the line segment between the current upper and lower bounds on marginal value. The bid table data is calculated according to optimality properties on resource assignments, and the constant elasticity of demand model implies higher revenue for lower prices and more connections provided. Therefore, the *best feasible* solution along the line segment between the upper and lower bounds on marginal value is where the largest number of connections are provided by assigning the maximum possible capacity. This in turn is obtained at the lowest feasible marginal values. Similarly, the *worst infeasible* solution along the line segment is found with the highest marginal values along that line, for which the capacity assignment is infeasible at every switch. We solve the following line search problems to obtain these two cases:

(P-feasible)

$$\text{Max } \alpha \quad (40)$$

subject to,

$$\sum_{i=1}^N \sum_{j=1}^{NS} \sum_{k=1}^{NS} BW_{ijk} (\bar{v}_1 - \alpha m_1, \bar{v}_2 - \alpha m_2, \dots, \bar{v}_{NS} - \alpha m_{NS}) IVP_{ijkx} < \overline{BW}_x \quad \forall x \quad (41)$$

(P-infeasible)

$$\text{Max } \beta \quad (42)$$

subject to,

$$\sum_{i=1}^N \sum_{j=1}^{NS} \sum_{k=1}^{NS} BW_{ijk} (v_1 + \beta m_1, v_2 + \beta m_2, \dots, v_{NS} + \beta m_{NS}) IVP_{ijkx} > \overline{BW}_x \quad \forall x \quad (43)$$

where,

$$m_j = \frac{\bar{v}_j - v_j}{\sum_{x=1}^{NS} \bar{v}_x - v_x} \quad (44)$$

$$0 < \alpha, \beta < \sum_{x=1}^{NS} \bar{v}_x - v_x \quad (45)$$

The line segment between the current set of lower bounds and upper bounds is given by (45). The maximum revenue feasible solution along this line is given by the set of marginal values  $\bar{v}_j - \alpha^* m_j$ . The minimum revenue infeasible solution is given by  $v_j + \beta^* m_j$ . We will use these solutions in the next step of the algorithm to determine which bounds should be changed before the next iteration of the algorithm.

**Step 3 Adjusting Lower and Upper Bounds on Marginal**

**Value:** Along the line segment between the upper and lower bounds, we wish to identify the switch for which the bandwidth which is most under-valued. That is, the switch(es) for which marginal value of bandwidth is “too low” at the solution  $\alpha^*$  to the revenue maximization problem (P-feasible) above, yielding the most assigned bandwidth. Similarly, we wish to select the switch (or switches) for which bandwidth is most over-valued, or the switch for which the lower bound is “too high” at the solution  $\beta^*$  to the problem (P-infeasible) above, yielding the least assigned bandwidth. These values are defined mathematically below:

$x_{\text{under-valued}}$

$$= \underset{x}{\operatorname{argmin}} \sum_{i=1}^N \sum_{j=1}^{NS} \sum_{k=1}^{NS} BW_{ijk} \left( \bar{v}_1 - \alpha^* m_1, \dots, \bar{v}_{NS} - \alpha^* m_{NS} \right) IVP_{ijkx} \quad (46)$$

$x_{\text{over-valued}}$

$$= \underset{x}{\operatorname{argmax}} \sum_{i=1}^N \sum_{j=1}^{NS} \sum_{k=1}^{NS} BW_{ijk} \left( v_1 + \beta^* m_1, \dots, v_{NS} + \beta^* m_{NS} \right) IVP_{ijkx} \quad (47)$$

The determination of which switches are over or under-valued in terms of their marginal values of bandwidth is the same as simply choosing the switch with the highest assigned bandwidth from the solution to (P-feasible) and that with the lowest assigned bandwidth to (P-infeasible). For example, if one switch is 100% assigned at the solution to (P-feasible) while all others are less than 50% assigned, the marginal value of bandwidth at the fully assigned switch is relatively too low for an optimal solution, since capacity should be fully assigned at all switches. Note that ties are permitted so that there may be more than one under-valued or over-valued switch.

To effectively raise the marginal value of bandwidth when solving (P-feasible), we will increase the lower bound on marginal value at that switch. When we then repeat the line search to solve (P-feasible) in the next iteration, a higher marginal value will result at that switch relative to the other switches. Similarly, we will decrease the upper bound on marginal value for the switch that has the lowest infeasible assignment at the optimal solution to (P-infeasible), e.g. lower the upper bound at a switch with 101% assignment when all others are ~150% assigned. This effectively lowers the marginal value calculated in the next iteration where we solve (P-infeasible) again.

$$\bar{v}_{x_{\text{over-valued}}} := \bar{v}_{x_{\text{over-valued}}} - \alpha^* m_{x_{\text{over-valued}}} \quad (48)$$

$$v_{x_{\text{under-valued}}} := v_{x_{\text{under-valued}}} + \beta^* m_{x_{\text{under-valued}}} \quad (49)$$

We raise the lower bound on marginal value for the under-valued switch to the value given by the solution to (P-infeasible), as given in (49). Similarly, we decrease the upper bound on marginal value for the over-valued switch by taking the value in the solution to (P-feasible), as given by (48). Note the use of an assignment operator, “:=” in (49) and (48). The value to the left of “:=” is the new value being calculated, while the value to the right of “:=” is based on the previous quantities.

Step 4 **Termination Test**: When the line search between the bounds yields a solution to (P-feasible) where capacity is fully assigned at all switches, we terminate the algorithm, and assign the values found in the line search problem (P-feasible) as the solution.

If,

$$\sum_{i=1}^N \sum_{j=1}^{NS} \sum_{k=1}^{NS} BW_{ijk} \begin{pmatrix} \bar{v}_1 - \alpha^* m_1, \dots \\ \bar{v}_{NS} - \alpha^* m_{NS} \end{pmatrix} IVP_{ijkx} \approx \overline{BW}_x \quad \forall x \quad (50)$$

then,

$$v_j^* = \bar{v}_j - \alpha^* m_j \quad \forall j. \quad (51)$$

Terminate the search.

Else,

Go to Step 2.

There are a couple of things to note about this iterative algorithm. First, the algorithm is really a *bisection search* as described in Step 3, where the gap between lower and upper bounds reduces at every step. It is well known that bisection searches are guaranteed to *converge*. We also show this numerically (See Figure 4 later). Next, because our search space is contained in the bid tables we have described earlier, which are indexed by integer variables,  $NCMAX_{ijk}$ , our search space is “lumpy”. As such, “fully assigned” means an arbitrarily chosen level of capacity utilization such as 99% assigned capacity at every switch. This is why we use the notation approximately equal, “ $\approx$ ”, rather than requiring strict equality, “ $=$ ”.

The final line search problem (P-infeasible) yields an upper bound on the objective value. That is for the local optimum given by (51), and the solution to (P-infeasible) we can calculate the expected revenues and state the duality gap for the local near-optimal solution is given as follows:

$$\sum_{i=1}^N \sum_{j=1}^{NS} \sum_{k=1}^{NS} \begin{pmatrix} REV_{ijk} (\underline{v}_1 + \beta^* m_1, \dots, \underline{v}_{NS} + \beta^* m_{NS}) \\ -REV_{ijk} (\bar{v}_1 + \alpha^* m_1, \dots, \bar{v}_{NS} + \alpha^* m_{NS}) \end{pmatrix} > 0 \quad (52)$$

## V. EXAMPLE PROBLEM

We will now present a simple real world problem with two-service classes and show how to set optimal prices and resource allocations using the model presented in section 3 along with the search method presented in section IV.

### A. Network Structure and Service Classes

Consider a single network offering voice and video connection to users via a set of five switches. We consider a bi-directional ring network (Fig. 3). The voice connections can be either local, i.e., routed through a single switch, or long-distance to anywhere in the network, i.e. routed from any origin to any destination switch. The video connections originate from a single switch (labeled “1” in the figure), where a video server is available and may be routed to any switch in the network, including the originating node. The connections are all routed using the minimum hop routing.

Switch capacities are chosen as follows. Switch “0” is included in the path for every video connection and has 5000

Mbps capacity, or five times as much capacity as switches “2” and “3”, which are each included in one fifth of the paths for

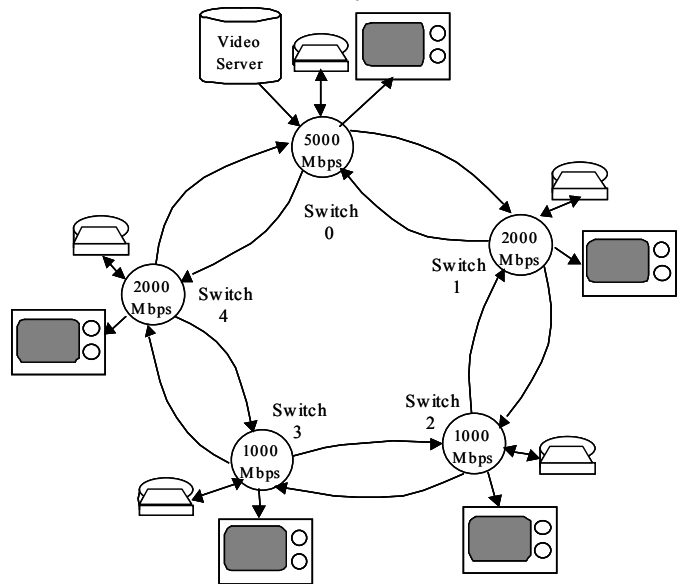


Fig. 3. Two service network with a single video server

TABLE II  
SERVICE CLASS DEFINITIONS

(a) Data Traffic Parameters				
	Mean. Rate (Mbps)	Peak. Rate (Mbps)	Average Burst (s)	
<b>Service Class</b>	$r_i$	$R_i$	$b_i$	
Voice	0.032	0.064	1.0	
Video	1.0	10.0	10.0	

(b) Other Parameters				
	Packet Loss Probability	Blocking Probability	Delay (s)	Average Holding Time (s)
<b>Service Class</b>	$PLOSS_i$	$BP_i$	$\bar{d}_i$	$T_i$
Voice	$10.0 \times 10^{-5}$	0.01	0.0	3.0
Video	$10.0 \times 10^{-9}$	0.01	5.0	30.0

video connections. Similarly, switches “1” and “4” have two times the capacity of “2” and “3”. The network services offered are voice and video. The service class definitions are given in Table II, which also provides traffic data

Voice service is a reflection of traditional voice service, which does not have a high peak bandwidth, and has relatively short periods of bursts. In terms of QoS, relatively high loss rates of packetized voice may be acceptable, but large delays cannot be tolerated. Voice connections are typically of short duration (e.g.  $T_{voice} = 3.0$  minutes), which results in a lower average number of connections in use for any arrival rate of requests, as defined by (3). Video is intended to reflect bursty sources of data traffic, where the bursts, such as action scenes can be a high peak and continue for prolonged periods. The peak and mean data rates are defined conservatively, based on MPEG-1 trace data in [29]. The length of the idle/busy period is chosen

based on a suggestion in [34] that video sessions see busy periods in the order of 10 seconds. For users with buffer space at the client end, a modest delay in the transmission of packets is acceptable. Some packet loss can be concealed, but a higher degree of reliability is required than for voice.

We assume constant elasticity of demand for both services, with the elasticity values taken from [18] (Table III). Note that, the elasticity of video is higher than that of voice. This means that the revenue, associated with video connections, increases more rapidly as the price is lowered than for voice connections. We have chosen the scaling factor in the functions to be 1 and 2 respectively. This means that at the price 1 (or a normalized price) the arrival rate of requests

TABLE III  
DEMAND FOR SERVICES

Service Class	Demand
Voice ( $\forall j,k$ )	$\lambda_{voice,jk} = 1.0p^{-1.05}$
Video ( $j = 1$ )	$\lambda_{video,jk} = 2.0p^{-2.50}$
Video ( $j \neq 1$ )	0

for video is twice as high, which reflects the higher value of a connection with video than that with simply voice.

### B. Illustration of Solution Algorithm

Fig. 4 and 5 provide details of the search steps when we implemented the algorithm for the Example Problem above. As we can observe, the search starts with lower and upper bound values at vectors  $\mathbf{0}$  and  $\mathbf{1}$  respectively, and through a number of bi-section steps, converges to the “best” (i.e. near optimal) solutions at iteration 8. All other variables also converge. To clarify exactly how the tight bounds are generated by the search method, we explain the step through an iteration in the Fig. 4 example, before continuing with the interpretation of the solution to the same example.

**A Single Iteration of the Search Algorithm:** At iteration “0”, Fig. 4 shows the starting upper and lower bounds on node marginal values. Recall from (32) that the marginal revenue per unit of bandwidth ( $u^+$ ) for a given service class must offset the sum of marginal value(s) of bandwidth ( $v$ ) at every along the path assigned to the service class. At each estimate of node marginal value, bandwidth assignments etc. will be obtained from the bid table. Thus, at the upper bounds on node marginal values, resource assignments will be very small and *feasible*. Conversely, at the lower bounds there will be *infeasible* over-assignment of resources. This is shown in Figure 5. At the initial upper bounds (1,1,1,1,1), bandwidth assignments at each node are (26.42%, 2.66%, 5.32%, 5.32%, 2.66%) of available capacity. At the lower bounds on node marginal values (0, 0, 0, 0, 0), the assignments are (899.32%, 988.73%, 1137.74%, 1137.74%, 988.73%) of available capacity. (These are not shown because they are outside the scale of Fig. 5).

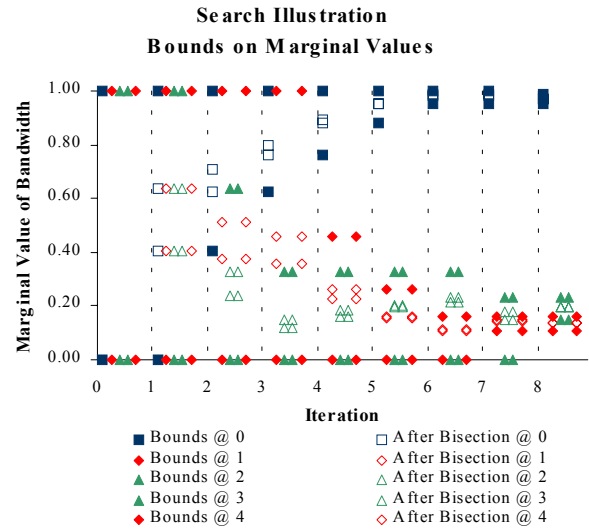


Fig. 4. Search Algorithm Steps: Convergence of Marginal Values

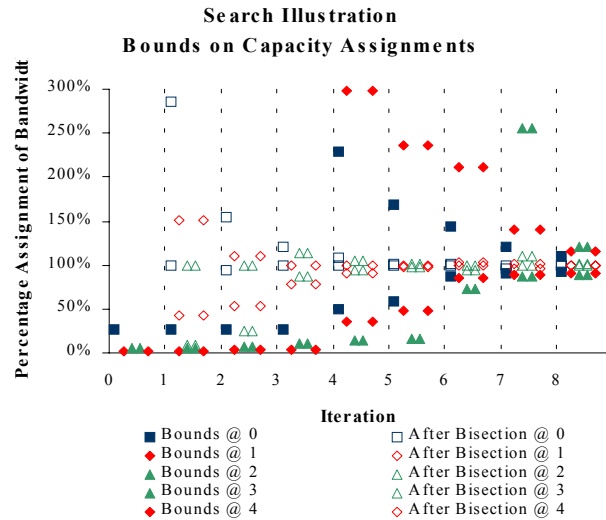


Fig. 5. Search Algorithm Steps: Convergence of Capacity Assignments

At iteration “1”, a *bisection* along the line segment joining (0, 0, 0, 0, 0) and (1, 1, 1, 1, 1) is obtained. If feasibility is not important, the maximum distance found along that line segment yields node marginal values to be (0.40, 0.40, 0.40, 0.40, 0.40); see Figure 4. In this case, the bandwidth assignments are (285.73%, 150.39%, 100.01%, 100.01%, 150.39%) of capacity (now shown in Figure 5). When feasibility has to be maintained, the bisection search yields node marginal values of (0.64, 0.64, 0.64, 0.64, 0.64), for which the bandwidth assignments are (99.98%, 43.04%, 8.10%, 8.10%, 43.04%); see Figure 5.

Given the broad range of capacity assignments (from 8.10% to 99.98 % of node capacity in the feasible case and from 100.01% to 285.73% of node capacity in the infeasible case), we have an intuition that the bounds on the node marginal value either under-assigns and over-assigns bandwidth at certain nodes in a relative sense. Node 0 with an infeasible assignment of 285.73% of capacity in this iteration is most

over-assigned, or what we now call, “under-valued”. Therefore, we raise the marginal value for node 0 relative to all the other nodes by resetting the lower bounds from (0, 0, 0, 0, 0) to (0.40, 0, 0, 0, 0). Likewise, the assignment at the tightest feasible solution is most under-assigned or “over-valued” at nodes 2 and 3, where only 8.10% of available capacity is included in the current feasible solution. Therefore, we lower the relative marginal value of nodes 2 and 3 by resetting the upper marginal value bounds from (1, 1, 1, 1, 1) to (1, 1, 0.64, 0.64, 1). The new bounds are passed to the algorithm for the next iteration and define the bounds at the beginning of iteration 2 in Fig. 4. The bisection search along the line segment between bounds on marginal values proceeds as above until convergence.

### C. Solution Interpretation

The solution is found when the algorithm converges to a set of marginal values on individual node bandwidth, which

TABLE IV  
OPTIMAL PRICE PER UNIT TIME FOR A CONNECTION OF SERVICE  $(i,j,k)$ .

Service Type	Origin ( <i>i</i> )	Destination ( <i>k</i> )				
		1	2	3	4	5
Voice	1	0.76	0.88	1.06	1.06	0.88
	2	0.88	0.10	0.24	0.40	1.00
	3	1.06	0.24	0.14	0.29	0.40
	4	1.06	0.40	0.29	0.14	0.24
	5	0.88	1.00	0.40	0.24	0.10
Video	1	1.65	1.91	2.32	2.32	1.91

occurs after 8 iterations in Fig. 4. The final marginal values of the nodes selected for the near optimal solution (i.e. the final solution to (P-feasible) as discussed in in Step 4 of Section IV are (0.97, 0.14, 0.20, 0.20, 0.14). The expected revenue at this solution is 11046.03. The duality gap (calculated by comparing the solution to (P-feasible) with the solution to (P-infeasible)) is roughly 0.24%; i.e. the best feasible solution we found was *within* 0.24% of the optimal solution. The bandwidth assignments at every switch are in excess of 99%, the criteria we used for termination in S4; i.e. the optimality conditions are, more or less, satisfied. Given the marginal values above, the optimal prices are simply looked up in the bid tables and are given in Table IV.

The prices reflect relative scarcity of bandwidth at switch “1”, due to the location of the video server. Relative to switch “1”, local (same origin and destination) voice connections are priced inexpensively elsewhere in the network. The long-distance connections must yield much higher marginal values than local connections to be profitable, since the long distance connections must outweigh the marginal value of local bandwidth all along the path of the connection. High marginal values correspond to higher prices and relatively smaller allocations. The video connections at switch “1” are the least expensive. Video connections at switch 1 use only resources at switch 1, while all other video connections must use this resource plus more resources at other switches. Consequently the volume of video connections at switch “1” must also be

the highest of anywhere in the network. However, the cheaper video service comes at the price of expensive voice service through this switch, because the voice connections over this switch must use resources with a very high marginal value relative to the other switches.

### D. Comparison to Flat Rate Pricing Approaches

To evaluate the usefulness of our approach, we compared optimal pricing with two flat rate pricing, per-hop and per-connection pricing, approaches. In the current US voice market, pricing per connection regardless of origin and destination is very popular. Pricing-per-hop is another flat-rate pricing scheme that has been proposed, and mimics the use of prices proportional to link use.

For both of these pricing approaches, we need only choose a price for each service type in the network, a total of  $N$  prices. Because we have only two service types, we assume a ratio between the voice price and video price and then simply solve for the smallest voice price that achieves feasibility at all switches. We will use the bid tables, with the resource allocations satisfying optimality properties and find the smallest feasible flat rate price for voice, subject to the assumed ratio between voice and video prices. Results are summarized in Figs. 6 and 7.

The optimal pricing mechanism provided greater revenue than either flat-rate pricing approach. The savings were upward of 20% in all cases. In general, pricing per hop produces higher revenue than pricing per connection in the reasonable ranges for relative video and voice prices. In this range, the expected revenue is relatively insensitive to the ratio between the prices for the two services. The pricing per connection shows the highest profit when voice and video connections are charged at the same rate. The prices essentially fill the network bandwidth with video connections, which are highly elastic and generate significant revenue as

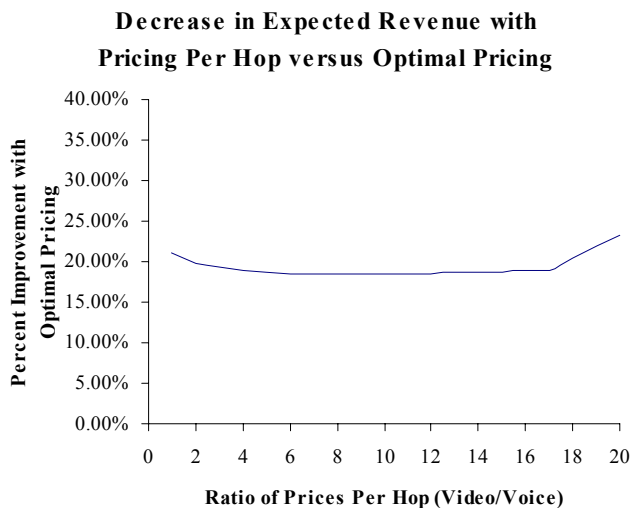


Fig. 6. Comparison of per hop pricing with optimal pricing.

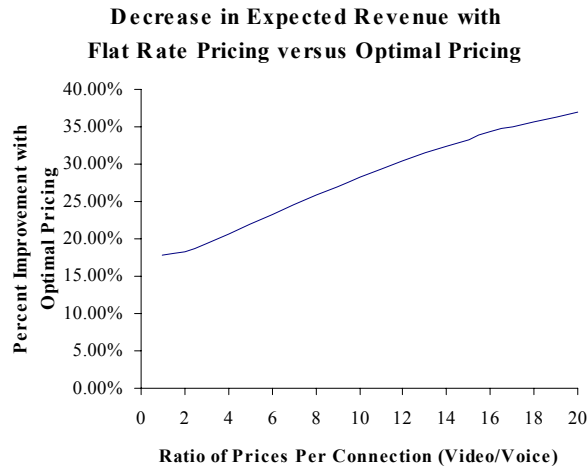


Fig. 7. Comparison of per connection pricing with optimal pricing.

prices are lowered. However, users are unlikely to accept a voice rate as high as would be necessary for a video connection. For this reason, the solutions with larger ratios between the video price and voice price are probably more practical. These solutions earn significantly less revenue.

## VI. CONCLUDING REMARKS

This paper presents a mathematical programming model for optimal pricing, and bandwidth and buffer allocation for multiple services with QoS guarantees in a connection-oriented network. We also present a novel solution methodology. We show, using numerical experiments, that flat-rate pricing (whether by connection or hop) is inferior to the multi-service pricing obtained from our model. Our model provides a very powerful mechanism for pricing multiple services in communications networks.

In order to implement our pricing scheme in a practical real world setting, the first decision to be made is whether one uses a centralized architecture or a decentralized one. In either case, a processor has to use demand (i.e. connection request) information to make projections of future demand for the different types of services. In a centralized architecture, the master switch will directly get information about capacity availabilities at each switch. Using this and the demand projections, the master switch will calculate node marginal values for each switch, and will determine the prices to be set and the resources to be allocated for each service class, updating it at reasonable time intervals. The model presented in this paper is particularly useful in this setting.

In a decentralized architecture, *each switch* will also have to calculate the prices and resources needed for each class of service. In order to have globally optimal solutions, each switch will have to communicate information pertaining to itself to the rest of the switches. Depending on whether the communication is sent only to neighboring switches, or to all of the switches, the actual calculations at any point in time might differ. The model we have presented in this paper will still work, if all of the decentralized switches use it for pricing

and resource allocation, but obtain perfect instantaneous information from the other switches on node marginal values and demand for bandwidth. However, if there is partial information and/or a significant time-lag for the communications protocol to yield up-to-date perfect information, a simultaneous implementation of the model by all switches might not yield a globally optimal solution.

We are currently working on multiple service pricing mechanisms for decentralized networks. We are also working on obtaining good demand projections, which is a further limitation of this paper that assumes that all parameters of the stochastic arrival process is known. The beginning of such work was reported in Keon's dissertation [15].

## APPENDIX: OPTIMALITY THEOREM

We argue for each of the optimality properties in the theorem individually:

i. Assume there is an optimal solution with prices and resource allocations such that the call-blocking constraint is non-binding, for at least one service class:

$$BP_{ijk}(\lambda_{ijk}(p_{ijk}), T_{ijk}, NCMAX_{ijk}) < \overline{BP}_{ijk} \quad (53)$$

It follows that for this class the network operator could allow a higher rate of connection requests and still satisfy the call-blocking constraint:

$$\exists \lambda'_{ijk}(p'_{ijk}) > \lambda_{ijk}(p_{ijk}) : BP_{ijk}(\lambda'_{ijk}(p'_{ijk}), T_{ijk}, NCMAX_{ijk}) \leq \overline{BP}_{ijk} \quad (54)$$

Demand is downward sloping, and marginal revenue with respect to prices is less than zero by assumption, so that the price corresponding to the higher rate of requests must be lower and revenue must be higher:

$$\frac{d\lambda_{ijk}(p_{ijk})}{dp_{ijk}} < 0, p'_{ijk} < p_{ijk} \quad (55)$$

$$BP_{ijk}(\lambda_{ijk}(p_{ijk}), T_{ijk}, NCMAX_{ijk}) = \overline{BP}_{ijk} \quad (56)$$

Therefore, for (53) an optimal solution cannot exist. The optimal solution must be such that call-blocking is binding:

$$BP_{ijk}(\lambda_{ijk}(p_{ijk}), T_{ijk}, NCMAX_{ijk}) = \overline{BP}_{ijk} \quad (24)$$

Assume there exists an optimal solution such that the allocation of bandwidth to a particular class is greater than the equivalent capacity for that class:

$$c(NCMAX_{ijk}, B_{ijk}, \overline{PLOSS}_{ijk}) < \frac{BW_{ijk}}{NCMAX_{ijk}} \quad (57)$$

There must exist a feasible allocation of less bandwidth, according to constraint (18) in (PNet) (the smallest feasible assignment is equal to equivalent capacity at the maximum permitted loss probability) for the particular class such that the QoS is still satisfied:

$$\exists BW'_{ijk} < BW_{ijk} :$$

$$c(NCMAX_{ijk}, B_{ijk}, \overline{PLOSS}_{ijk}) = \frac{BW'_{ijk}}{NCMAX_{ijk}} \quad (58)$$

The reassignment given by (58), where equivalent capacity,  $c_i$ , is minimized by selecting the loss probability equal to its

maximum permitted value, may make it possible to assign bandwidth for an additional connection, lower the price of such a class and increase revenue. Even if admitting more connections is not possible due to the reassignment in (58), revenue cannot decrease from the reassignment. Therefore, the assumption that bandwidth is excess of effective bandwidth, (57), makes no sense and the optimal solution is invalid. We must have assigned bandwidth equal to equivalent capacity:

$$c(NCMAx_{ijk}, B_{ijk}, PLOSS_{ijk}) = \frac{BW_{ijk}}{NCMAx_{ijk}} \quad (25)$$

ii. Assume that a solution is optimal and possesses the following property for at least one local service class, e.g.  $(x,j)$ :

$$\left( \begin{aligned} & (NCMAx_{xjj} + 1)c(NCMAx_{xjj} + 1, B_{xjj}, \overline{PLOSS}_{xjj}) \\ & + \sum_{i=1}^N \sum_{k=1}^{NS} NCMAx_{ijk} c_{ijk}(NCMAx_{ijk}, B_{ijk}, \overline{PLOSS}_{ijk}) \end{aligned} \right) \leq \overline{BW}_j \quad (59)$$

For the particular service class  $(x,j)$ , we can allocate resources for an additional connection, and tolerate a higher arrival rate of requests, at a lower price, while satisfying all constraints:

$$\exists \lambda'_{xjj}(p'_{xjj}) > \lambda_{xjj}(p_{xjj}), p'_{xjj} < p_{xjj} : \text{all constraints satisfied} \quad (60)$$

Similarly to above, a lower price results in increased revenue. Therefore (59) is false and capacity must be fully allocated:

$$\left( \begin{aligned} & (NCMAx_{xjk} + 1)c_{xjk}(NCMAx_{xjk} + 1, B_{xjk}, \overline{PLOSS}_{xjk}) \\ & + \sum_{i=1}^N \sum_{k=1}^{NS} NCMAx_{ijk} c_{ijk}(NCMAx_{ijk}, B_{ijk}, \overline{PLOSS}_{ijk}) \end{aligned} \right) > \overline{BW}_j \quad (26)$$

$1 \leq x \leq N, 1 \leq j \leq NS$

We first assume that there exists an optimal solution that satisfies the following property for a single service class  $(i,j,k)$ , which contradicts condition (27) in Theorem 1:

$$B_{ijk} \bar{d}_{ijk} < BW_{ijk} \quad (61)$$

For this service class, as we increase the buffer allocation,  $B_{ijk}$ , the equivalent capacity  $c$ , or minimum bandwidth assignment per connection falls with the increased buffer space:

$$\frac{\partial c(B_{ijk})}{\partial B_{ijk}} < 0 \quad (62)$$

Therefore, there exists buffer and bandwidth allocations different from those given in (61), for which the bandwidth assignment is lower:

$$\exists B'_{ijk} > B_{ijk}, BW'_j = c_j(B'_j) < BW_j : B'_{ijk} = \bar{d}'_{ijk} BW'_j \quad (63)$$

If the difference in bandwidth is sufficiently large there may be a revenue increasing solution, by (26) above, meaning the current solution cannot be optimal. On the other hand if the bandwidth allocation does not yield an improvement in revenue, the revenue cannot be decreased by the reallocation described above. Therefore, the optimal solution must satisfy

(27):

$$B_{ijk} \bar{d}_{ijk} = BW_{ijk}, \quad 1 \leq i \leq N \quad (27)$$

#### ACKNOWLEDGMENT

We thank without implicating Roch Guerin, Yannis Korilis and Vijay Chandru and referees of this journal for comments and suggestions

#### REFERENCES

- [1] Aldebert, M., M. Ivaldi and C. Roucolle, "Telecommunications Demand and Pricing Structure: An Econometric Analysis," *Proceedings of the 7th International Conference on Telecommunications Systems: Modelling and Analysis*, pp. 255-267, Nashville, TN, March 18-21, 1999.
- [2] Bertsekas, D. and R. Gallager, *Data Networks*, Prentice-Hall, Englewood Cliffs, New Jersey, (1987).
- [3] Bertsekas, D., "Auction Algorithms for Network Flow Problems: A Tutorial Introduction," *Computational Organization and Applications*, vol. 1, 7-66 (1992).
- [4] Bhat, U.N. and Basawa, I.V., *Queuing and Related Models*, Oxford University Press, New York, New York (1992).
- [5] Cocchi, R., S. Shenker, D. Estrin and L. Zhang, "Pricing in Computer Networks: Motivation, Formulation, and Example," *IEEE/ACM Transactions on Networking*, vol. 1, no. 6, 614-627 (1993).
- [6] Courbetis, C., V. Siris and G.D. Stamoulis, "Integration of Pricing and Flow Control for Available Bit Rate Services in ATM Networks," *Globecom*, (1996).
- [7] Economides, A.A. and J.A. Silvester, "Multi-Objective Routing in Integrated Services Networks: A Game Theory Approach," in *Proceedings of INFOCOM'91*, pp. 1220-1225 (1991).
- [8] Edell, R.J., N.McKeown and P.P. Varaiya, "Billing Users and Pricing for TCP," *IEEE Journal on Selected Areas in Communications*, vol. 13, pp. 1162-1175, September 1995.
- [9] Guérin, R., H. Ahmadi and M. Naghshineh, "Equivalent Capacity and its Application to Bandwidth Allocation in High-Speed Networks," *IEEE Journal on Selected Areas in Communications*, vol. 9, no. 7, 968-981 (1991).
- [10] Gibbens, R.J. and F.P. Kelly, "Resource Pricing and the Evolution of Congestion Control," draft paper (1999). Available: <http://www.statslab.cam.ac.uk/~frank/evol.html>.
- [11] Jiang, H. and S. Jordan, "The Role of Price in the Connection Establishment Process," *European Transactions on Telecommunications - Economics of Telecommunications*, November 9 (1994).
- [12] Jordan, S. and H. Jiang, "A Pricing Model for High Speed Networks with Guaranteed Quality of Service," *IEEE InfoCom*, pp. 888-895, March 1996.
- [13] Kelly, F.P., "Tariffs and Effective Bandwidths in Multiservice Networks," *Proceedings of ITC'94*, pp. 401-410, 1994.
- [14] Kelly, F.P. A.K. Maullo and D.K.H. Tan, "Rate Control in Communication Networks: Shadow Prices, Proportional Fairness and Stability," *Journal of the Operational Research Society*, vol. 49, 237-252 (1998). Available: <http://www.statslab.cam.ac.uk/~frank/rate.html>.
- [15] Keon, Neil J., "Pricing in Telecommunications Networks Offering Multiple Services and Quality of Service Guarantees", University of Pennsylvania, Systems Engineering Department, August 2000.
- [16] Korilis, Y.A. and A.A. Lazar, "On the Existence of Equilibria in Noncooperative Optimal Flow Control," *Journal of the ACM*, vol. 42, pp. 584-613, May 1995.
- [17] Korilis, Y.A., T.A. Varvarigou and S.R. Ahuja, "Pricing Noncooperative Networks," submitted to the *IEEE/ACM Transactions on Networking*, May 1997.
- [18] Lanning, S., D.Mitra, Q. Wang and M. Wright, "Optimal Planning for Optical Transport Networks," presented at the *Fifth INFORMS Telecommunications Conference*, Boca Raton, FL, March 5-8, 2000.
- [19] Low, S.H. and P. Varaiya, "A New Approach to Service Provisioning in ATM Networks," *IEEE/ACM Transactions on Networking*, vol. 1, no. 5, 547-553 (1993).

- [20] Low, S.H. and D. E. Lapsley, "Optimization Flow Control, I: Basic Algorithm and Convergence," *IEEE/ACM Transactions on Networking*, Dec. 1999.
- [21] Low, S. H., "Equilibrium Bandwidth and Buffer Allocations for Elastic Traffics," *IEEE/ACM Transactions on Networking*, Vol. 8, no. 3, pp 373-383, (2000)
- [22] Mackie-Mason, J.K. and H. Varian, "Pricing the Internet," *Public Access to the Internet*, eds. B. Kahin and J. Keller, Cambridge and London: MIT Press, 269-314, 1995.
- [23] Mackie-Mason, J.K. and H.R. Varian, "Pricing Congestible Network Resources," *IEEE Journal on Selected Areas in Communications*, vol. 13, pp. 1141-1149, September 1995.
- [24] Mackie-Mason, J.K. and H. Varian, "Some Economics of the Internet," in: *Networks, Infrastructure and the New Task for Regulation*, eds. W. Sichel and D.L. Alexander, Ann Arbor: University of Michigan Press, 107-36, (1996).
- [25] Masuda, Y. and S. Whang, "Dynamic Pricing for Network Service: Equilibrium and Stability," *Management Science*, vol. 45, no. 6, pp. 857-869, (1999).
- [26] Mendelson, H. and S. Whang, "Optimal Incentive-Compatible Priority Pricing for the M/M/1 Queue," *Operations Research*, vol. 38, no. 5, pp. 870-883, (1990).
- [27] Paschalidis, I.Ch., and John N. Tsitsiklis, "Congestion-Dependent Pricing of Network Services," *IEEE/ACM Transactions on Networking*, vol. 8, no. 2, pp 171-184 (2000)
- [28] Perros, H.G. and K.M. Elsayed, "Call Admission Control Schemes: A Review," *IEEE Communications*, vol. 34, no. 11, 82-91 (1996).
- [29] Rose, O., "Statistical Properties of MPEG Video Traffic and Their Impact on Traffic Modeling in ATM Systems", *Proceedings of the 20th Annual Conference on Local Computer Networks*, Minneapolis, MN, pp 397-406, 1995.
- [30] Shenker, S., "Service Models and Pricing Policies for an Integrated Services Internet," *Public Access to the Internet*, eds. B. Kahin and J. Keller, Cambridge and London: MIT Press, 315-337 (1995).
- [31] Thomas, P., Teneketzis, D. and Mackie-Mason, J.K., "A Market-Based Approach to Optimal Resource Allocation in Integrated-Services Connection-Oriented Networks," submitted for publication, August 1999.
- [32] de Veciana, G., G. Kesidis and J. Walrand, "Resource Management in Wide-Area ATM Networks Using Effective Bandwidths," *IEEE Journal on Selected Areas in Communications*, vol. 13, no. 6, 1081-1089 (1995).
- [33] de Veciana, G. and R. Baldick, "Resource Allocation in Multi-Service Networks via Pricing: Statistical Multiplexing," *Computer Networks and ISDN Systems*, vol. 30, 951-962 (1998).
- [34] Walrand, J. and P. Varaiya, *High-Performance Communication Networks*, Morgan Kaufmann Publishers, Inc. San Francisco, California (1996).
- [35] Wang, Q., J.M. Peha, and M.A. Sirbu, "The Design of an Optimal Pricing Scheme for ATM Integrated Services Networks," *Special Issue: Internet Economics, Journal of Electronic Publishing*, University of Michigan Press, vol. 2, no. 1 (1996).

He is an Associate Editor of *Operations Research, Management Science, Telecommunications Systems, Network and Spatial Economics*, and *Computers and Operations Research*. He is a Senior Member of INFORMS (Institute For Operations Research and Management Science), and on the Advisory Boards of a number of Information Technology start-up companies.

**Neil J. Keon** became a Member (M) of IEEE in 2000. Neil has a Bachelor of Com. with joint honours in Economics and Finance (1992) from McGill University, an M.Sc in Systems Science. (1995) from the University of Ottawa, and an M.S.E (1999) and a Ph.D. (2000) in Systems Engineering from the University of Pennsylvania.

He is an Assistant Professor in Information Technology and Operations Management at the Cox School of Business, Southern Methodist University, Dallas TX, where he has been since 2000.

He is also a member of INFORMS.

**G. Anandalingam** (M'85, SM'90) has a B.A. (1975) from Cambridge University in Electrical Sciences, and Ph.D. (1981) in operations research with economics from Harvard University.

Anand is the *Ralph J. Tysler Professor of Management Science* at the Robert H. Smith School of Business, and the Institute for Systems Research at the University of Maryland where he has been since 2001. From 1987-2001, he was the *National Center Professor* (Engineering School), and *Professor of Operations & Information Management* (Wharton School) at the University of Pennsylvania. He has published more than 70 papers on telecommunications and networking issues, and a book on Telecommunications Network Design to be published by Kluwer Scientific Publishers.