

Improving Utilization for Deterministic Service In Multimedia Communication

Hui Zhang
Lawrence Berkeley Laboratory
MailStop: 50B-229
Berkeley, CA 94720
hzhang@george.lbl.gov

Domenico Ferrari
University of California at Berkeley
and
International Computer Science Institute
ferrari@icsi.berkeley.edu

Abstract

Network-based real-time multimedia applications require guaranteed performance communication services. To provide guaranteed service, resources have to be reserved within the network. If reservation is based on the peak rate of each connection, the network will be under-utilized by guaranteed service traffic when the traffic is bursty. In this paper, we first show that local deterministic delay bounds can be guaranteed over a link for bursty traffic even when the sum of the peak rates of all the connections is greater than the link speed. Compared to previous admission control conditions, the new result allows a multi-fold increase in the number of admitted connections when the traffic is bursty. We then show that this new result can be efficiently extended from a single switch to a network of arbitrary topology by using rate-controlled service disciplines at the switches.

1 Introduction

High speed networking has introduced opportunities for new network-based multimedia applications such as video conferencing, scientific visualization and medical imaging. These applications have stringent performance requirements in terms of throughput, delay, delay jitter and loss rate. The best-effort service provided by the current packet-switching networks is not adequate. New services are needed.

Two types of new services have been proposed in the literature to support real-time multimedia applications: *guaranteed service* [9] and *predicted service* [3]. In a guaranteed service model, client-specified *a priori* performance bounds are guaranteed to each connection *regardless of* the behaviors of other connections. In a predicted service model, a network dictated *post facto* delay bound and playback point are provided to the communication client. The playback point may vary and the service may be disrupted due to the network load fluctuation. It is assumed

that applications using the predicted service can adapt to the changing of the playback point and tolerate infrequent service disruptions.

There are two important factors to consider in choosing between these two service models: the quality of the service and the cost of providing such a service. From a communication client's point of view, the quality of the guaranteed service is better since the performance bounds are guaranteed and there are no service disruptions. From the network's point of view, the cost of providing a service depends on the maximum network utilization that can be achieved. Although it is argued that the predicted service would allow a higher network utilization by real-time traffic than guaranteed service [3], since there is no complete description of the algorithms for the predicted service, especially the admission control algorithms, there are no quantitative comparisons available.

In [6], two types of guaranteed services are proposed: *deterministic service* and *statistical service*. In deterministic service, performance bounds are guaranteed for *all* packets on a connection even in the worst case. In statistical service, *probabilistic* performance bounds are guaranteed. Although the quality of deterministic service is better, statistical service allows the network to achieve a higher utilization by exploiting statistical multiplexing.

It is clear from the above discussion that achieving high network utilization is one of the most important considerations in providing guaranteed service, especially deterministic service. To provide guaranteed service, resources have to be reserved within the network. In general, resource reservation schemes can achieve high network utilizations for smooth traffic. However, many clients requiring performance guarantees have bursty sources, for example, compressed video. If reservation is based on the peak rate of each connection, new requests will be rejected when the sum of the peak rates of all the connections reaches link speed. In this case, the network will be under-utilized by guaranteed service traffic when the peak-to-average-rate

ratio is high.

Therefore, it is important to derive solutions that can provide performance guarantees even when the sum of the peak rates of all the connections is greater than the link speed, and to understand the relationship between the server utilization and traffic burstiness. In this paper, we discuss these issues within the framework of the Tenet real-time channel scheme. We show that deterministic guarantees can be provided even when the sum of peak rates of all the connections is greater than the link speed, and reasonable average network utilization can be achieved for deterministic service even when traffic is bursty. Compared to previous admission control algorithms for deterministic service, the new result allows a multi-fold increase in the number of admitted connections for a single server when the traffic is bursty.

Providing local performance bounds at a single switch only solves part of the problem. In a networking environment, packets from different connections are multiplexed at each of the switches. Even if the traffic can be characterized at the entrance to the network, complex interactions among connections will destroy many properties of the traffic inside the network, and the traffic model at the source may not be applicable inside the network. Since local performance bounds can be guaranteed for a connection only if the connection’s input traffic to the switch satisfies certain traffic characterization, traffic pattern distortion may make it difficult to guarantee local performance bounds at switches inside the network. We address this problem by using rate-controlled service disciplines inside the network to reconstruct traffic patterns. Rate-controlled service disciplines allows the result for a single switch be efficiently extended to a network of arbitrary topology.

The remainder of the paper is organized as follows. Section 2 briefly reviews the real-time channel scheme and shows the limitations of previous admission control tests for the deterministic service. Section 3 presents the analysis which provides delay bounds for a First-Come-First-Served (FCFS) scheduler. Section 4 gives numerical examples to illustrate the results derived in Section 3. Section 5 discusses issues on providing end-to-end delay bounds in a networking environment. Section 6, discusses the implications of the results. Finally, Section 7 concludes the paper with a summary.

2 Background

In this section, we give a brief overview of the current version of the Tenet resource management algorithms [9, 8].

The Tenet algorithms are based on a communication abstraction called a *real-time channel* [9]. A real-time channel is a network connection associated with traffic and performance parameters. The parameters are provided by

$\frac{Smax}{Xmin \times l}$	peak rate
$\frac{Smax}{Xave \times l}$	upper bound on average rate over l
$\frac{Xave}{Xmin}$	peak-to-average-rate ratio or burst ratio
$\frac{l}{Xave} Xmin$	maximum burst length

Table 1: ($Xmin, Xave, l, Smax$) Traffic Model

the clients to specify their traffic characteristics and performance requirements. The traffic specification consists of parameters ($Xmin, Xave, l, Smax$), where $Xmin$ is the minimum packet inter-arrival time, $Xave$ is the worst-case average packet inter-arrival time over an averaging interval, l is the averaging interval, and $Smax$ is the maximum packet size. In such a specification, it is easy to see that $Xave/Xmin$ is the peak-to-average-rate ratio, which is an indicator of the traffic burstiness. It should be also noticed that l also affects the traffic burstiness. For the given values of $Xmin$ and $Xave$, l determines how long the source can continuously send packets at the peak rate in the worst case. The larger $Xave/Xmin$ and l , the burstier the traffic. Table 1 shows the model and the interpretations of some of the formulas. In the table, l denotes the link speed.

Tenet algorithms provide both deterministic and statistical guarantees. For the purpose of this paper, we only consider deterministic guarantees. For a channel with deterministic delay bound \bar{D} , the network guarantees that delays of all packets on that channel will be less than \bar{D} as long as the channel does not violate its traffic specification.

A channel needs to be *established* before data can be transmitted. This channel establishment is achieved in the following manner: a real-time client specifies its traffic characteristics and end-to-end performance requirements to the network; the network determines the most suitable route for the channel, translates the end-to-end parameters into local parameters at each node, and attempts to reserve resources at these nodes accordingly. This is done in a distributed manner during a round-trip communication.

In order to provide performance guarantees, two levels of controls are needed: at the channel level, channel admission control algorithms reserve resources for each of the channels and limit the maximum utilization of the network by real-time traffic; at the packet level, the service discipline at each of the switches determines the multiplexing policy and allocates resources to different channels according to the reservations.

As shown in [7, 21], many service disciplines can be used to provide real-time service. However, different service disciplines require different admission control algorithms. In [9], conditions are given for a variation of the Earliest-Due-Date discipline. Three tests need be satisfied before a new channel request can be accepted, which are: (a) *deterministic test*, which limits the peak utilization of the

server to be less than 1; (b) *delay bound or schedulability test*, which avoids the case of scheduling saturation; and (c) *buffer test*, which ensures enough buffer space is available. For the purpose of the discussion in this paper, we will focus on the deterministic test¹, which requires

$$\sum_{j=1}^n \frac{1}{Xmin_j} \times \frac{Smax_j}{l} < 1 \quad (1)$$

where l is the link speed and n is the number of channels traversing the link.

Although, the tests guarantee that deterministic delay bounds can be provided, they are rather restrictive in the sense that the sum of the peak rates of all deterministically guaranteed connections on any link has to be less than the link speed. If the peak-to-average-rate ratio is high, the average link utilization by real-time traffic will be low.

3 Delay Analysis

In this section, we show that deterministic delay bounds can be obtained even when the sum of peak rates of all real-time channels is greater than the link speed, though the sum of the average rates of all real-time channels has to be less than the link speed. The result holds even for a simple discipline like FCFS. We use the bounding techniques developed by Cruz [4]. In [4], a fluid traffic model (σ, ρ) is used. A channel satisfies traffic specification (σ, ρ) if, for any time interval of length u , the number of bits arriving during the interval is no more than $\sigma + \rho u$. The model we use is the discrete $(Xmin, Xave, I, Smax)$ model.

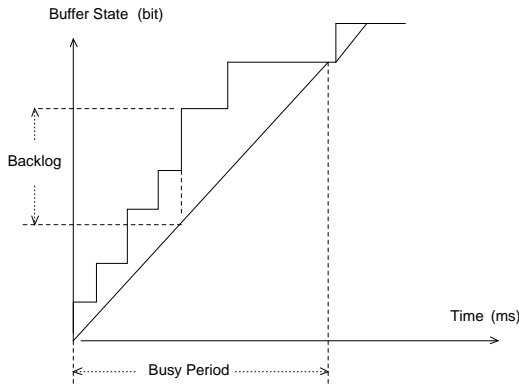


Figure 1: Concepts: Delay, Backlog and Busy Period

Figure 1 illustrates some of the concepts used in the analysis. The horizontal axis is the time. The upper curve in the figure is the sum of the bits that have arrived since the origin of time. Each arrival of a packet causes an upward

¹ We consider a network with non-blocking switches, where queueing happens at the output link of each switch. The processing time of each packet is thus equal to the transmission time, and is bounded by $\frac{Smax_j}{l}$

jump in this curve. The lower curve is the number of bits that have been transmitted. The difference between the two curves is the *backlog* function. A work-conserving server always transmits packets at a constant rate when there is a backlog. The points where the arrival curve and the service curve meet, or where the backlog is zero, divide the time axis into *busy periods* and *idle periods*. Within such a framework, the following two propositions immediately follow, where \overline{Smax} is the maximum packet size that can be transmitted over the link, and l is the link speed.

Proposition 1 For a work-conserving real-time scheduler, if for any realization of the input traffic that satisfies a given traffic constraint, the maximum length of a busy period is no greater than \overline{d} , then $\overline{d} + \overline{Smax}/l$ is an upper bound for the delays of all packets.

Proposition 2 For a FCFS real-time scheduler, if for any realization of the input traffic that satisfies a given traffic specification, the maximum real-time traffic backlog divided by the link speed is no greater than \overline{d} , then $\overline{d} + \overline{Smax}/l$ is an upper bound for the delays of all packets.

To take into account of the effect of non-real-time packets, which have a lower priority than real-time packets, but cannot be preempted after the beginning of their transmission, \overline{Smax}/l is included in the delay bounds.

Notice that in both propositions, in order to derive the delay bound, certain upper bound values (the length of busy period or the backlog) need to be calculated for any realization of the input traffic. It is impossible to calculate these upper bounds for every realization of the input traffic because there are infinite such realizations. If we can find one “worst-case” realization of the input traffic, and prove that this realization gives the maximum upper bound values, we can just analyze the scheduler for this realization and compute the delay bound.

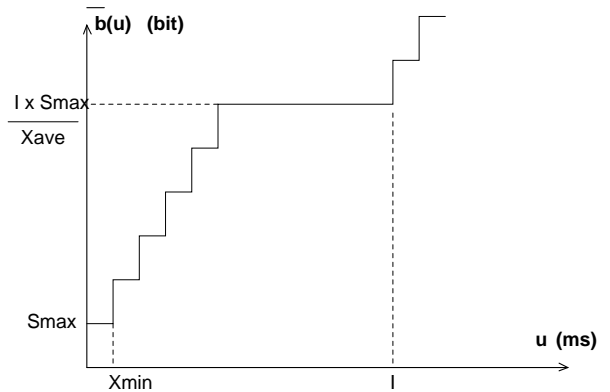


Figure 2: Traffic Constraint Function for $(Xmin, Xave, I, Smax)$ Specification

Central to the analysis is the concept of *traffic constraint function* $\bar{b}(\cdot)$ ². $\bar{b}_j(u)$ is defined to be the maximum number of bits that can arrive on channel j during *any* interval of length u . For a traffic source that obeys the model $(Xmin_j, Xave_j, I_j, Smax_j)$, it is easy to show that $\bar{b}_j(u)$ is

$$(\min(\lceil \frac{u \bmod I_j}{Xmin_j} \rceil, \lceil \frac{I_j}{Xave_j} \rceil) + \lceil \frac{u}{I_j} \rceil \lceil \frac{I_j}{Xave_j} \rceil) Smax_j$$

Figure 2 illustrates the function $\bar{b}(\cdot)$ for a channel with traffic specification $(Xmin, Xave, I, Smax)$.

The following theorem bounds the delay for a FCFS scheduler.

Theorem 1 *Let there be n channels multiplexed on a link with a FCFS scheduler and link speed l . If for $j = 1, \dots, n$, the traffic on channel j is bounded by $\bar{b}_j(\cdot)$, then the delays of packets on all the channels are bounded by \bar{d} , where \bar{d} is defined by*

$$\bar{d} = \frac{1}{l} \max_{u \geq 0} \left\{ \sum_{j=1}^n \bar{b}_j(u) - l \times u + \frac{\overline{Smax}}{l} \right\} \quad (2)$$

Proof.

Let ξ be a realization of the input traffic. If s is the starting time instant of a busy period, and t is a time instant within the busy period, define $b_j^\xi(s, t)$ to be the number of bits that arrive during the interval (s, t) on channel j . The maximum backlog in this realization ξ is $\max_{s, t} \{ \sum_{j=1}^n b_j^\xi(s, t) - l \times (t - s) \}$. We have

$$\begin{aligned} & \frac{1}{l} \max_{s, t} \{ \sum_{j=1}^n b_j^\xi(s, t) - l \times (t - s) \} \\ & \leq \frac{1}{l} \max_{s, t} \{ \sum_{j=1}^n \bar{b}_j(t - s) - l \times (t - s) \} \\ & \leq \frac{1}{l} \max_{u \geq 0} \{ \sum_{j=1}^n \bar{b}_j(u) - l \times u \} \end{aligned}$$

Since ξ is an arbitrary realization of the input traffic, from Proposition 2. \bar{d} as defined by Equation (2) is a delay bound for packets from all channels. **Q.E.D.**

In the following two corollaries, we give closed-form solutions for two special cases. Corollary 1 considers the case when $\sum_{j=1}^n \frac{Smax_j}{Xmin_j} \leq l$, Corollary 2 considers the case of homogeneous sources when $\sum_{j=1}^n \frac{Smax_j}{Xmin_j} > l$.

Corollary 1 *Let there be n channels multiplexed on a link with a FCFS scheduler and link speed l . Assume channel j obeys the traffic specification $(Xmin_j, Xave_j, I_j, Smax_j)$, ($j = 1, \dots, n$). If $\sum_{j=1}^n \frac{Smax_j}{Xmin_j} \leq l$, then the delays of packets on all the*

²We use a notation that differs slightly from that in [4], where $b(\cdot)$ is used to denote the traffic constraint function. In this paper, $\bar{b}(\cdot)$ is used to denote the traffic constraint function, and $b^\xi(s, t)$ denotes the number of bits that arrive between time instants s and t for an input traffic realization ξ .

channels are bounded by \bar{d} , where \bar{d} is given by

$$\bar{d} = \frac{1}{l} \sum_{j=1}^n Smax_j + \frac{\overline{Smax}}{l} \quad (3)$$

Proof.

$$\begin{aligned} \bar{d} &= \frac{1}{l} \max_{u \geq 0} \left\{ \sum_{j=1}^n \bar{b}_j(u) - l \times u \right\} + \frac{\overline{Smax}}{l} \\ &\leq \frac{1}{l} \max_{u \geq 0} \left\{ \sum_{j=1}^n \lceil \frac{u}{Xmin_j} \rceil Smax_j - l \times u \right\} + \frac{\overline{Smax}}{l} \\ &\leq \frac{1}{l} \max_{u \geq 0} \left\{ \sum_{j=1}^n \left(\frac{u}{Xmin_j} + 1 \right) Smax_j - l \times u \right\} + \frac{\overline{Smax}}{l} \\ &= \frac{1}{l} \sum_{j=1}^n Smax_j + \frac{\overline{Smax}}{l} + \left(\sum_{j=1}^n \frac{Smax_j}{Xmin_j} - l \right) \frac{u}{l} \\ &\leq \frac{1}{l} \sum_{j=1}^n Smax_j + \frac{\overline{Smax}}{l} \end{aligned}$$

Q.E.D.

In the proof, the first inequality holds due to the assumption that the minimum inter-packet spacing is $Xmin_j$. All three inequalities become equalities when $u = 0$, i.e., the worst case backlog is at the beginning of a busy period, when every channel sends out a maximum length packet.

Corollary 2 *Let there be n homogeneous channels with traffic specification $(Xmin, Xave, I, Smax)$ multiplexed on a link with the FCFS service discipline and link speed l . If $n \times \frac{Smax}{Xave} \leq l$ and $n \times \frac{Smax}{Xmin} > l$, the delays of packets from all the channels are bounded by \bar{d} , where \bar{d} is given by*

$$d = Xmin + \left(\mu_{ave} - \frac{1}{burstRatio} \right) I + \frac{\overline{Smax}}{l} \quad (4)$$

where $\mu_{ave} = n \times Smax / Xave \times l$, and $burstRatio = Xave / Xmin$.

Proof.

Let \hat{I} be the length of the longest interval that packets from one channel can arrive at the scheduler with $Xmin$ spacing, we have,

$$\hat{I} = \min\{u \mid \bar{b}(u) = \bar{b}(I)\} \quad (5)$$

$$= \left(\frac{I}{Xave} - 1 \right) Xmin \quad (6)$$

It is easy to see that $\sum_{j=1}^n (\bar{b}_j(u) - l \times u)$ is maximized at

$u = \hat{I}$.

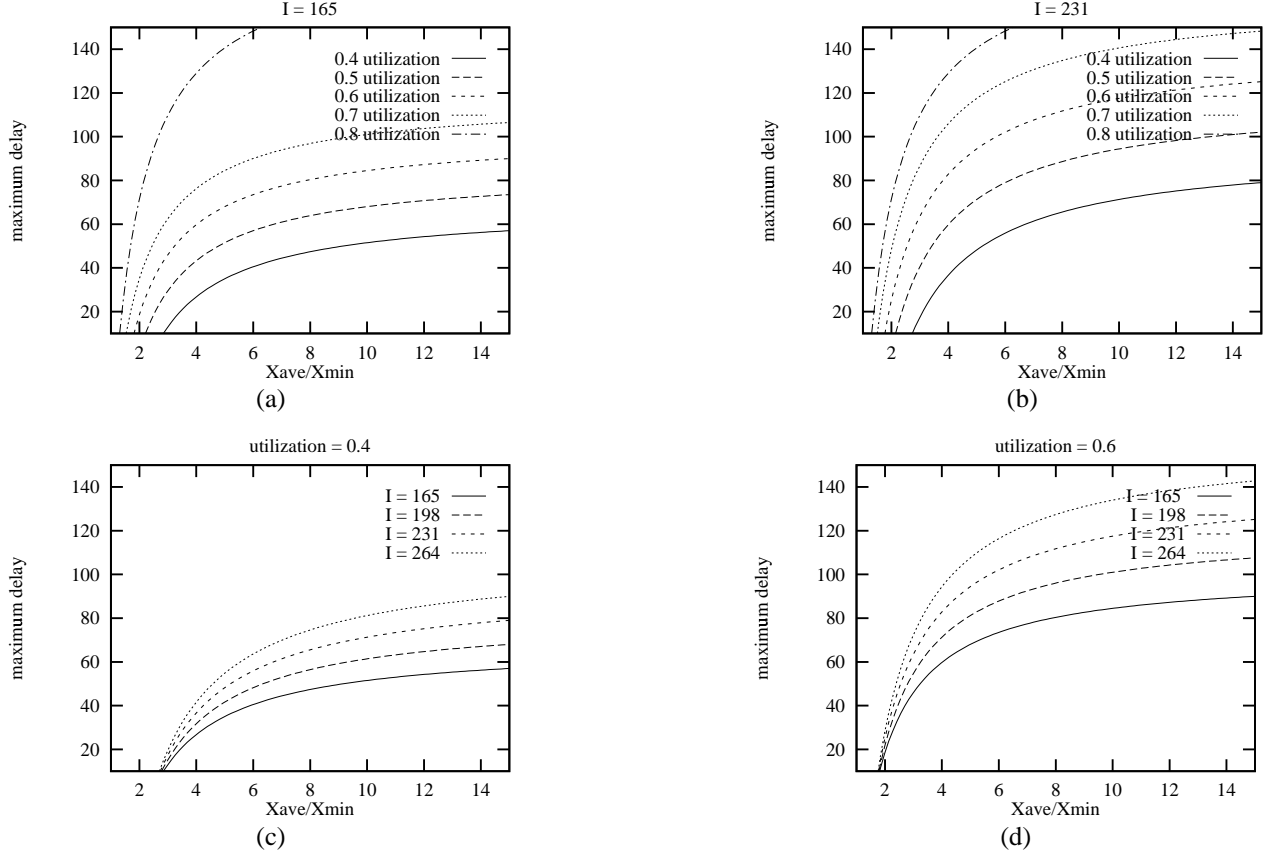


Figure 3: Effects of Burstiness and Utilization on Delay Bound

From Theorem 1, we have,

$$\begin{aligned}
\bar{d} &= \frac{1}{l} \sum_{j=1}^n (b_j(\hat{I}) - l \times \hat{I}) + \frac{\overline{Smax}}{l} \\
&= \frac{1}{l} \left(\sum_{j=1}^n Smax \frac{I}{Xave} \right) - \hat{I} + \frac{\overline{Smax}}{l} \\
&= \mu_{ave} \times I - \hat{I} + \frac{\overline{Smax}}{l} \\
&= \mu_{ave} \times I - \left(\frac{I}{Xave} - 1 \right) Xmin + \frac{\overline{Smax}}{l} \\
&= Xmin + \left(\mu_{ave} - \frac{1}{burstRatio} \right) I + \frac{\overline{Smax}}{l}
\end{aligned}$$

Q.E.D.

In (4), μ_{ave} is the average utilization of the link, and $Xave/Xmin$ and I represent the burstiness of the channels. Intuitively, the more loaded the link, and the burstier the channel traffic, the larger the delay bound. This is exactly what is shown by (4). It is important to see that the averaging interval I affects the traffic burstiness. For the given values of $Xmin$ and $Xave$, I determines how long the source can continuously send packets at the peak rate in the worst case. This is illustrated in Equation (6).

The previous two corollaries consider only two special cases. More general result is given by the following corollary:

Corollary 3 Let there be n channels multiplexed on a link with a FCFS scheduler and link speed l . Assume channel j obeys the traffic specification $(Xmin_j, Xave_j, I_j, Smax_j)$, ($j = 1, \dots, n$). If $\sum_{j=1}^n \frac{Smax_j}{Xave_j} \leq l$, then the delays of packets on all the channels are bounded by \bar{d} , where \bar{d} is given by

$$\bar{d} = \frac{\overline{Smax} + \sum_{j=1}^n \frac{Smax_j}{Xave_j} [I_j (1 - \frac{Xmin_j}{Xave_j}) + Xmin_j]}{(1 - \mu_{ave})l}$$

$$\text{where } \mu_{ave} = \sum_{j=1}^n \frac{Smax_j}{Xave_j \times l}.$$

The detail of the proof is given in [19]. Though Corollary 3 is more general than Corollary 1 and 2, the bounds given by Corollary 1 and 2 in those special cases are tighter than those given by Corollary 3.

4 Numerical Examples

In this section, we show some numerical examples to illustrate the results presented in the previous section. For

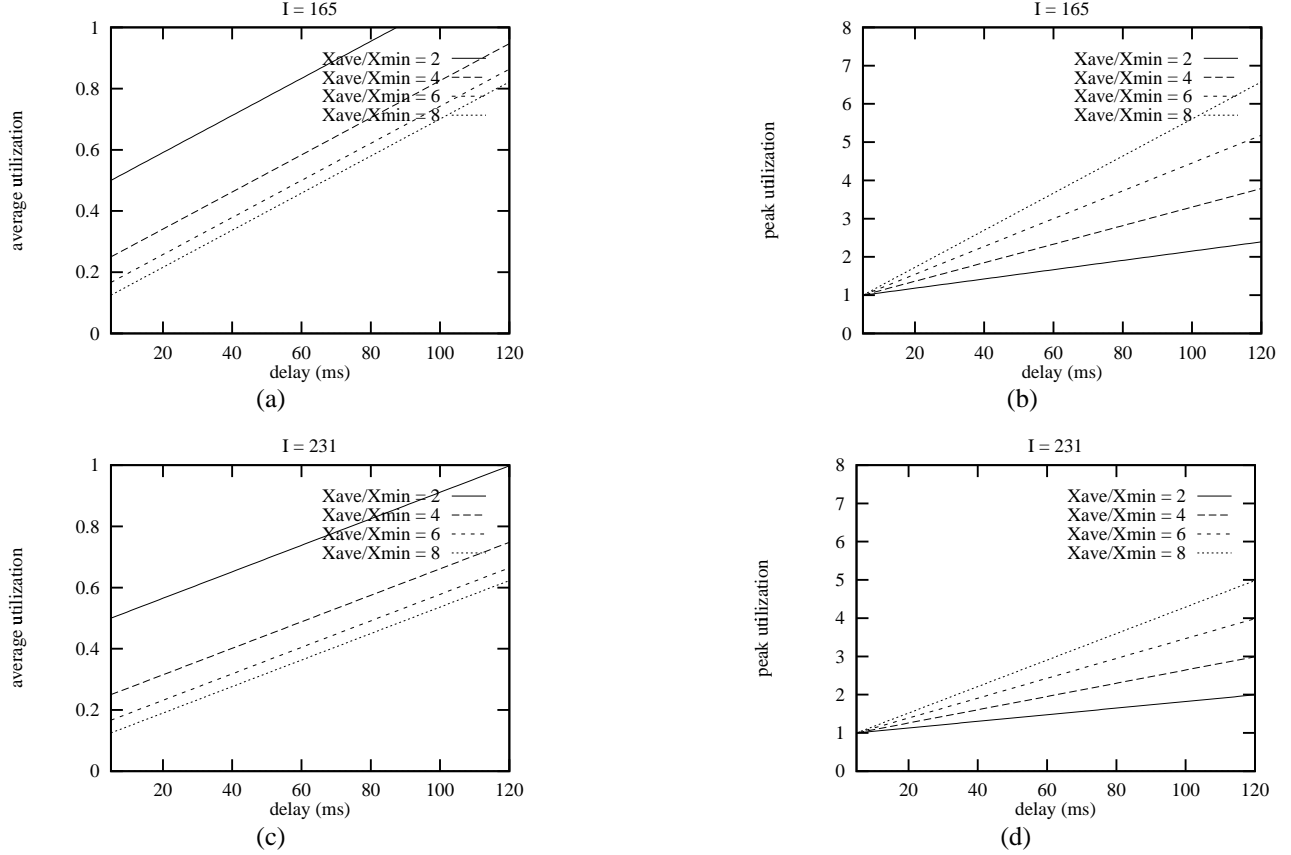


Figure 4: Average and Peak Utilization vs. Delay Bound

simplicity, we consider only the situation where all the sources have the same traffic characteristics. For heterogeneous sources, Corollary 3 can be used to derive delay bounds.

Figure 3 illustrates the effects of the traffic burstiness and the average utilization on delay bounds. The horizontal axis is the peak-to-average-rate ratio of the channels, and the vertical axis is the delay bound. The figures only show cases when the sum of the peak rates of all the channels is greater than 1. The first thing to be noticed from these figures is that, when the average utilization μ_{ave} and the averaging interval I are fixed, the delay bound that can be provided increases as the peak-to-average-rate ratio increases, i.e., burstier traffic results in a larger delay bound. A dual result is that for the same delay bound, a higher average utilization can be achieved when the traffic is less burstier. This is shown in Figure 4. Figures 3(a) and 3(b) show the results with fixed averaging interval I but different average utilizations. It can be seen that, when the averaging interval and the peak-to-average-rate ratio are fixed, a higher average utilization of the link results in a larger delay bound, i.e., the more loaded the link, the larger the delay bound. Figures 3(c) and 3(d) show the results for fixed utilization but different I 's. The values of I are

chosen to be multiples of 33 ms, which is the time interval between two frames in a 30 frames/sec video stream. It can be seen that, when the average utilization and the peak-to-average-rate ratio are fixed, a larger averaging interval results in a larger delay bound. The intuition is that with the same peak-to-average-rate ratio, a larger averaging interval means a burstier traffic; under the same average utilization of the link, a burstier traffic results in a larger delay bound.

Figure 4 shows how much average or peak utilization can be achieved under certain delay bound constraints, where average and peak utilization are defined by $\sum_{j=1}^n \frac{S_{max_j}}{X_{ave_j} \times I}$ and $\sum_{j=1}^n \frac{S_{max_j}}{X_{min_j} \times I}$, respectively. Here *average utilization* means the maximum fraction of bandwidth that can be allocated to real-time traffic. The actual bandwidth used by real-time traffic is below the allocated bandwidth. The bandwidth unused by real-time traffic can be used by non-real-time traffic. Since the maximum peak utilization of the link is 1 under the old deterministic test defined in Equation (1), the peak utilization can be seen as an improvement factor of how many more channels can be accepted under the new admission control algorithms than under the ones with the deterministic test.

As can be seen, the peak utilization of the link by deterministically guaranteed performance traffic can be much

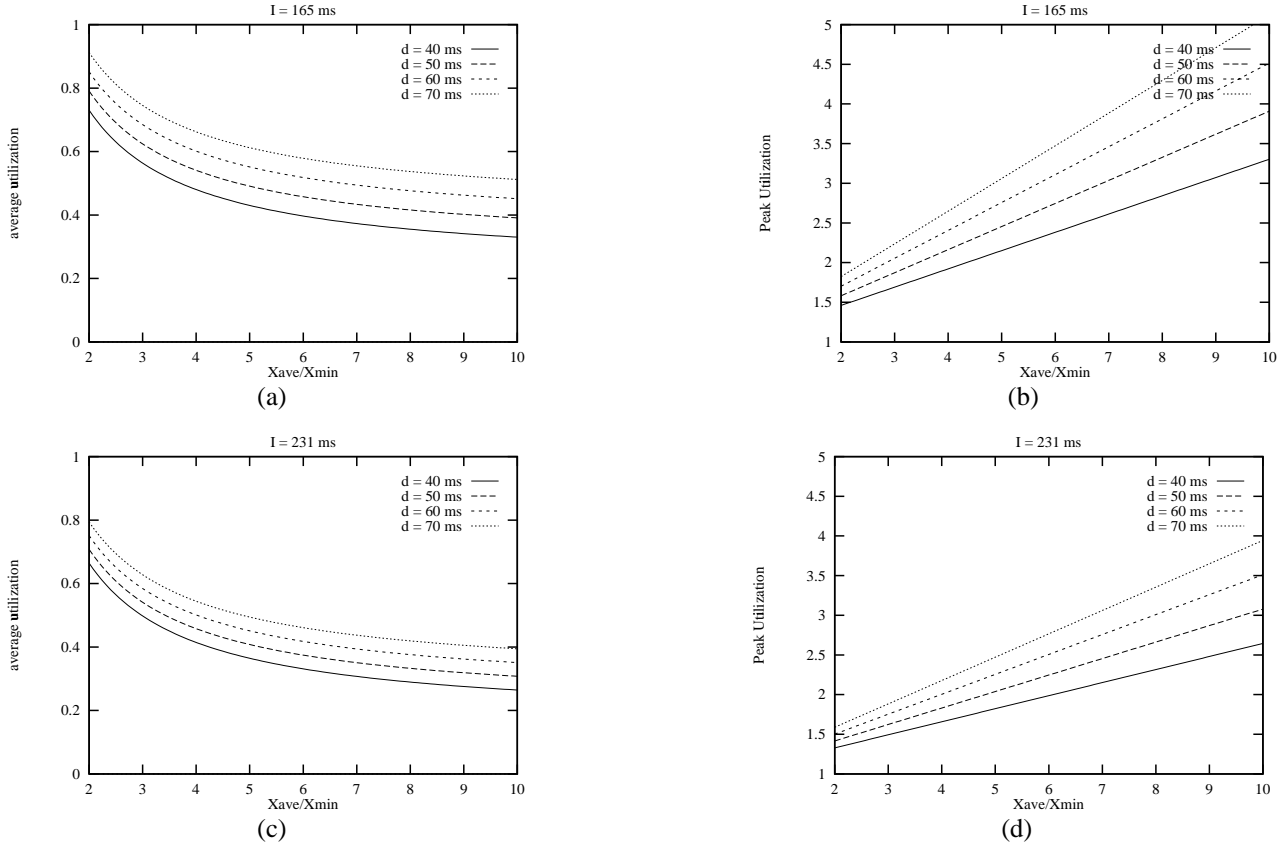


Figure 5: Average and Peak Utilization vs. Burst Ratio

greater than 1. Compared to the admission criteria where the sum of the peak rates of all channels have to be no greater than the link speed, using the admission control condition as stated in Theorem 1 may provide a multi-fold increase in the number of admitted channels when the traffic is bursty.

Figure 5 shows how average and peak utilization are affected by the burst ratio. As can be seen, on the one hand, the peak utilization, or the statistical multiplexing gain, is higher when the burst ratio is higher; on the other hand, the average utilization is lower when the burst ratio is higher. This matches our intuition that higher burst ratios provide more opportunity for statistical multiplexing but result in lower network utilization in order to meet a specific performance guarantee. Similar results can be obtained for statistical service [19, 22].

5 Providing end-to-end deterministic guarantees

In the previous section, we showed that even for simple service disciplines like FCFS, deterministic delay bounds can be obtained when the sum of the peak rates of all the channels is greater than the link bandwidth. However, the result holds only for a single scheduler. In this section, we extend the analysis from a single scheduler to a network of

switches, and show that end-to-end deterministic guarantees can be provided in general networking environments.

In a networking environment, packets from different channels are multiplexed at each switch. Even if the traffic can be characterized at the entrance to the network, complex interactions among channels will distort the traffic pattern and destroy many properties of the traffic inside the network. Thus, the traffic model at the source may not be applicable inside the network. Since local performance bounds can be guaranteed for a channel only if the channel's input traffic to the switch satisfies certain traffic characterization, traffic pattern distortion may make it difficult to guarantee local performance bounds at switches inside the network.

One solution to this problem is to characterize the traffic pattern distortion inside the network, and derive the traffic characterization at the entrance to each switch from the characterization of the source traffic and the traffic pattern distortion. This approach, taken in [4, 1, 15, 14], has several limitations.

First, it only applies to networks with *constant* delay links. Constant delay links have the desirable property that the traffic pattern at the receiving end of the link is the same as that at the transmitting end of the link. This property is important for these solutions because central to the

analysis is the technique of characterizing the output traffic from a scheduler and using it as the input traffic to the next-hop scheduler. However, in an internetworking environment, links connecting switches may be subnetworks such as ATM or FDDI networks. Though it is possible to bound delay over these subnetworks, the delays for different packets will be *variable*. Thus, these solutions do not apply to an internetworking environment.

Second, most of these solutions apply only to a restricted class of networks. Characterizing the traffic pattern inside the network is equivalent to solving a set of multi-variable equations [5, 15, 14]. In a feedback network, where traffic from different channels form traffic loops, the resulting set of equations may be unsolvable. Thus, most of these solutions apply only to feed-forward networks or a restricted class of feedback networks.

Finally, in networks with *work-conserving* service disciplines, even in the situations when traffic inside the network can be characterized, the traffic is usually more bursty inside the network than that at the entrance. This is independent of the traffic model being used. In [4], a deterministic fluid model (σ, ρ) is used to characterize traffic source. A source is said to satisfy (σ, ρ) if during any time interval of length u , the amount of its output traffic is less than $\sigma + \rho u$. In such a model, σ is the maximum burst size, and ρ is the average rate. If the traffic of channel j is characterized by (σ_j, ρ_j) at the entrance to the network, its characterization will be $(\sigma_j + \sum_{i'=1}^{i-1} \rho_j \bar{d}_{i',j}, \rho_j)$ at the entrance to the i -th switch along the path, where $\bar{d}_{i',j}$ is the local delay for the channel at the i' -th switch. Compared to the characterization of the source traffic, the maximum burst size at the i -th switch increases by $\sum_{i'=1}^{i-1} \rho_j \bar{d}_{i',j}$. This increase of burst size grows linearly along the path.

In [14], a family of stochastic random variables is used to characterize the source. Channel j is said to satisfy a characterization of $\{(R_{t_1,j}, t_1), (R_{t_2,j}, t_2), \dots\}$, where $R_{t_i,j}$ are random variables, and $t_1 < t_2 < \dots$ are time intervals, if $R_{t_i,j}$ is *stochastically larger* than the number of packets generated over any interval of length t_i by source j . If the traffic channel j is characterized by $\{(R_{t_1,j}, t_1), (R_{t_2,j}, t_2), \dots\}$, at the entrance to the network, its characterization will be $\{(R_{t_1 + \sum_{i'=1}^{i-1} b_{i',j}}, t_1), (R_{t_2 + \sum_{i'=1}^{i-1} b_{i',j}}, t_2), \dots\}$ at the i' -th switch, where $b_{i'}$ is the maximum busy period at switch i' . The same random variable that bounds the maximum number of packets over an interval at the entrance of the network, now bounds the maximum number of packets over a much *smaller* interval at switch j . I.e., the traffic is *burstier* at switch j than at the entrance.

In both the (σ_j, ρ_j) and $\{(R_{t_1,j}, t_1), (R_{t_2,j}, t_2), \dots\}$, analysis, the burstiness of a channel's traffic accumulates at each hop along the path from the source to destination, which results in a low utilization of the network.

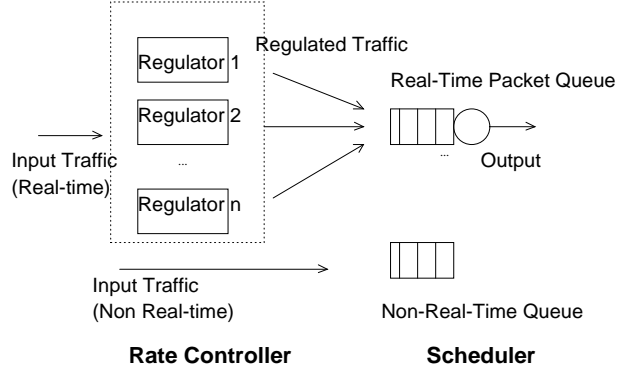


Figure 6: Rate-Controlled Service Disciplines

Another solution to the traffic pattern distortion problem, which we adopt in our approach, is to reconstruct the traffic pattern at each switch with a class of non-work-conserving service disciplines called rate-controlled service disciplines [19]. As shown in Figure 6, a rate-controlled service discipline consists of two components, a rate-controller and a scheduler. The rate-controller shapes the input traffic from each channel into the desired traffic pattern by assigning an eligibility time to each packet. The scheduler then orders the transmission of eligible real-time packets from all channels. Many types of regulators and schedulers can be used. Different combinations of regulators and schedulers result in different service disciplines. The class is quite general. Most non-work-conserving disciplines proposed for high speed networks such as Delay-EDD [18], Stop-and-Go [11], Hierarchical Round Robin [12], and Rate-Controlled Static Priority [20], either belong to this class, or can be implemented by a rate-controlled service discipline with the appropriate choices of rate-controllers and schedulers [19].

Rate-controlled service disciplines have the following two properties:

- (1) If a channel's traffic satisfies certain traffic characteristics at the entrance to the network, with use of the appropriate rate-controllers, the same characteristics will be satisfied by the traffic at the entrance to each of schedulers along the path. This allows us to perform delay analysis at each scheduler using the *same* traffic characterization.
- (2) The end-to-end delay of a packet in a network with rate-controlled servers consists of the following components: waiting time in the schedulers, holding time in the rate-controllers and the link delays. In [19], it is shown that the end-to-end delay can be bounded by the sum of bounds on link delays and bounds on waiting time in the schedulers; holding packets in rate-controllers will not increase the end-to-end delay bound, although it may increase the end-to-end average delay.

Properties (1) and (2) are significant. Property (1) means that we can analyze the delay characteristic of each scheduler along a path with the *same* traffic characteristics of the original source. The traffic characteristics need not be $(X_{min}, X_{ave}, I, S_{max})$ discussed in this paper. For example, if a channel can be characterized by a MMPP at the entrance to the network, it can be characterized by the *same* MMPP at each of the schedulers. Property (2) means that we can combine the delay analysis of each individual scheduler and obtain the end-to-end delay characteristics of a channel. This applies to both deterministic and statistical analysis [19].

6 Discussion

In previous sections, we showed that even for simple service disciplines like FCFS, deterministic delay bounds can be obtained when the sum of the peak rates of all the channels is greater than the link bandwidth. The question naturally arises: why does one need more complex service disciplines than FCFS? There are several reasons that FCFS alone is not enough.

First, guaranteed service requires that the network *protect* clients from two sources of variability: misbehaving users and network load fluctuations. The FCFS discipline does not offer such protections. Although putting traffic policing function at *all* network access points will prevent misbehaving users from affecting other users, traffic distortions due to network fluctuations suggest that protection should be implemented within the network using rate-based service disciplines [21], for example, rate-controlled disciplines as discussed in Section 5.

The second reason FCFS is not enough is that an FCFS server can only offer a single value of delay bound for all the channels. However, the performance requirements for integrated services networks will be diverse. It is important to support multiple classes of Quality of Service. If only one delay bound is provided as in the case of FCFS, the delay bound has to satisfy the most stringent requirement among all the channels, and will therefore under-utilize the network when the requirements of channels are different. In [20, 19], it has been shown that the Rate-Controlled Static Priority (RCSP), which consists of a rate-controller and a static priority scheduler, strikes a good balance between simplicity of implementation and flexibility in allocating bandwidths and delay bounds to different channels; also, it can achieve a reasonably high average utilization for deterministic real-time traffic even when channels have different performance requirements.

The third reason FCFS may not be enough is that more sophisticated service disciplines can provide better bounds than the FCFS discipline.

In the paper, we have shown that reasonably high average utilization can be achieved for deterministic service

even when the traffic is bursty. This is only true when the peak-to-average-rate ratio and the averaging interval are relatively small. In our examples, the values we chose for peak-to-average-rate ratio were 2 to 8, and the values we chose for the averaging interval were 99 ms to 298 ms. It should be noticed that these numbers are reasonable for compressed video. Recent video traffic trace study shows that the peak-to-average-rate ratio for VBR video is 1.5 - 4 [2, 13]. Also, for MPEG [10], the interval between two intra-frame-coding frames (*I* frames) is normally between 3 to 9 frame sizes, which corresponds to 99 ms to 298 ms when the video is played at 30 frames per second.

However, there are situations where providing only deterministic service may significantly underutilize the network. Two solutions can be adopted to enhance the utilization of the network by guaranteed performance traffic. If applications can tolerate certain losses of data without significantly affecting the quality, *statistical services* can be provided to achieve a higher average network utilization by exploiting statistical multiplexing [9, 17, 19]. Also, *cooperative, consenting, high-level multiplexing* schemes can be used to address the tradeoffs between the quality of services offered to each individual clients and the overall utilization of the network [16].

7 Conclusion

In this paper, we have showed that it is possible to provide deterministic delay bounds even when the sum of the peak rates of all the channels is greater than the link speed. Even for a simple discipline like FCFS, a reasonable average utilization can be achieved for deterministic service. Compared to the previous deterministic test, the new result allows a multi-fold increase in the number of connections that can be accepted when the traffic is bursty. We have showed that the improvement factor increases as the burst ratio becomes higher; however, the overall average link utilization is lower when the burst ratio is higher. By using rate-controlled service disciplines, we efficiently extend the result to general networking environment of arbitrary topology, which includes both feedback and feed-forward networks, and internetworks with variable but bounded link delays.

8 Acknowledgements

The authors are grateful to Jorg Liebeherr for helpful discussions at the beginning of this research.

References

- [1] A. Banerjea and S. Keshav. Queueing delays in rate controlled networks. In *Proceedings of IEEE INFOCOM'93*, pages 547–556, San Francisco, CA, April 1993.

- [2] A. W. Berger, S. P. Morgan, and A. R. Reibman. Statistical multiplexing of layered video streams over ATM networks with leaky-bucket traffic descriptors, 1993. preprint.
- [3] D. Clark, S. Shenker, and L. Zhang. Supporting real-time applications in an integrated services packet network: Architecture and mechanism. In *Proceedings of ACM SIGCOMM'92*, pages 14–26, Baltimore, Maryland, August 1992.
- [4] R. L. Cruz. A calculus for network delay, part I : Network elements in isolation. *IEEE Transaction of Information Theory*, 37(1):114–121, 1991.
- [5] R.L. Cruz. A calculus for network delay, part II : Network analysis. *IEEE Transaction of Information Theory*, 37(1):121–141, 1991.
- [6] D. Ferrari. Client requirements for real-time communication services. *IEEE Communications Magazine*, 28(11):65–72, November 1990.
- [7] D. Ferrari. Real-time communication in an internet-work. *Journal of High Speed Networks*, 1(1):79–103, 1992.
- [8] D. Ferrari, A. Banerjee, and H. Zhang. Network support for multimedia: a discussion of the Tenet approach. Technical Report TR-92-072, International Computer Science Institute, Berkeley, California, October 1992. Also to appear in *Computer Networks and ISDN Systems*.
- [9] D. Ferrari and D. Verma. A scheme for real-time channel establishment in wide-area networks. *IEEE Journal on Selected Areas in Communications*, 8(3):368–379, April 1990.
- [10] D. Le Gall. MPEG: A video compression standard for multimedia applications. *Communications of the ACM*, 34(4):46–58, April 1991.
- [11] S. J. Golestani. A stop-and-go queueing framework for congestion management. In *Proceedings of ACM SIGCOMM'90*, pages 8–18, Philadelphia Pennsylvania, September 1990.
- [12] C.R. Kalmanek, H. Kanakia, and S. Keshav. Rate controlled servers for very high-speed networks. In *IEEE Global Telecommunications Conference*, pages 300.3.1 – 300.3.9, San Diego, California, December 1990.
- [13] H. Kanakia, P. Mishra, and A. Reibman. An adaptive congestion control scheme for real-time packet video transport. In *Proceedings of ACM SIGCOMM'93*, pages 20–31, San Francisco, California, September 1993.
- [14] J. Kurose. On computing per-session performance bounds in high-speed multi-hop computer networks. In *ACM SigMetrics'92*, 1992.
- [15] A.K.J. Parekh and R.G. Gallager. A generalized processor sharing approach to flow control in integrated services networks: The multiple node case. In *Proceedings of the INFOCOM'93*, pages 521–530, San Francisco, CA, March 1993.
- [16] C. Parris, G. Ventre, and H. Zhang. Graceful adaptation of guaranteed performance service connections. In *Proceedings of IEEE GLOBECOM'93*, Houston, TX, November 1993.
- [17] D. Verma. *Guaranteed Performance Communication in High Speed Networks*. PhD dissertation, University of California at Berkeley, November 1991.
- [18] D. Verma, H. Zhang, and D. Ferrari. Guaranteeing delay jitter bounds in packet switching networks. In *Proceedings of Tricomm'91*, pages 35–46, Chapel Hill, North Carolina, April 1991.
- [19] H. Zhang. Service disciplines for integrated services packet-switching networks. PhD Dissertation. UCB/CSD-94-788, University of California at Berkeley, November 1993.
- [20] H. Zhang and D. Ferrari. Rate-controlled static priority queueing. In *Proceedings of IEEE INFOCOM'93*, pages 227–236, San Francisco, California, April 1993.
- [21] H. Zhang and S. Keshav. Comparison of rate-based service disciplines. In *Proceedings of ACM SIGCOMM'91*, pages 113–122, Zurich, Switzerland, September 1991.
- [22] H. Zhang and E. Knightly. Providing end-to-end statistical performance guarantees with interval dependent stochastic models. In *ACM Sigmetrics'94*, Nashville, TN, May 1994.