

BUFFER SPACE ALLOCATION FOR REAL-TIME CHANNELS IN A PACKET-SWITCHING NETWORK

Domenico Ferrari and Dinesh C. Verma

Computer Science Division

Department of Electrical Engineering and Computer Sciences

University of California

and International Computer Science Institute

Berkeley, California, U.S.A.

Abstract

Broadband integrated networks will have to offer *real-time communication services*; that is, they will have to transport information with performance guarantees. A paper previously published by the authors presented a scheme for establishing real-time channels in a pure packet-switching network; that scheme did not include any method for allocating buffer space in the network's nodes to the channels being established. This paper completes the description and evaluation of that scheme, since it presents one such method, and some of the results of the extensive simulations performed to test it. The method is found to be correct and to have a low overhead. While the utilization of the buffer space allocated to the statistical channels is often quite low, thereby indicating that our worst-case approach tends to overallocate space to those channels, the space our method gives to deterministic channels seems to be reasonably well utilized.

This research has been supported in part by AT&T Bell Laboratories, Hitachi, Ltd., the University of California under a MICRO grant, the National Science Foundation under Grant No. CDA-8722788, and the International Computer Science Institute. The views and conclusions in this document are those of the authors and should not be interpreted as representing official policies, either expressed or implied, of any of the sponsoring organizations.

1. Introduction

Broadband integrated networks will have to offer *real-time communication services*, i.e., to transport information with performance guarantees, primarily to satisfy the requirements of digital audio and digital video [Lein89]. The specific performance guarantees that will be needed will vary with the type of traffic (video, audio, or their various combinations with each other and with other information media) as well as with the application [Ferr90]. Once such services become available, they are likely to be used even for types of traffic that do not include video or audio, for example, alarm conditions in distributed real-time systems, urgent remote database queries, or file transfers requiring lower bounds on throughput. Thus, real-time communication services are going to be essential in the integrated networks of the future, along with the more traditional, non-real-time ones such as datagram and virtual circuit services of various kinds.

Three basic approaches can be used to provide real-time communication: circuit switching, packet switching, and hybrid switching. A *circuit-switching* network will allocate its resources continuously and exclusively to each connection for the entire duration of the connection [Harr80]. With *packet switching*, resource allocation is neither continuous nor exclusive; that is, resources are shared among various conversations on a packet-by-packet basis, and the only privileges certain conversations may receive consist of a higher priority (preemptive, but with a non-preemptable granule, i.e., the packet) in the use of shared resources. In a *hybrid-switching* network, the two previously described approaches coexist: certain conversations (typically, those needing performance guarantees) are carried by the circuit-switching portions, others by the packet-switching portion of the network.

If a sufficient amount of resources is allocated to a connection in a circuit-switching network, the given performance guarantees can be satisfied. However, bursty traffic will make a poor use of these dedicated resources. A hybrid solution could solve this problem if real-time traffic were non-bursty, and if bursty traffic did not need performance guarantees; this, unfortunately, is not the case, since compressed video and audio are bursty, and since other bursty traffic is likely to require real-time services. It is therefore useful to investigate the question whether packet switching can be used to provide performance guarantees, and, if so, whether and under what conditions a packet-switching solution is more convenient than the other two.

In [Ferr89] and [FeVe90], we have provided an answer to the question of feasibility by proposing and evaluating a scheme on which the design of a real-time service for an integrated wide-area packet-switching network could be based. One of the assumptions made in that paper was that the amount of buffer space available in each node would be sufficient to hold all the packets waiting to be forwarded to the next node, even in the worst case. In this paper, we eliminate that assumption, and describe and evaluate an algorithm for determining the amount of buffer space that should be allocated to each real-time connection in each node on the basis of the characteristics of the traffic on that connection and the bound on the probability of packet loss specified by the client. The algorithm completes the scheme described in [Ferr89] and [FeVe90], and the constructive proof of feasibility that scheme was intended to provide.

The important features of the scheme (under the infinite buffer space assumption) are summarized in Section 2. Section 3 presents the buffer space allocation algorithm. The simulation experiments we have run to evaluate the algorithm and their results are discussed in Section 4. Finally, Section 5 provides a summary of our study and suggestions for future work.

Related work can be found in [CoYa88], [Zhan89] and [OhON88].

2. A Scheme for Providing Real-Time Services in a Packet-Switching Network

The solution proposed in [Ferr89] and [FeVe90] is based on the notion of *real-time channel*, a simplex connection over which packets can travel from the sender to the receiver with performance satisfying given guarantees. Our scheme allows the following performance guarantees to be offered to the clients of the service:

- (1) *throughput bounds*: both
 - (i) minimum peak throughput (θ_{\min}) to be provided by the network, and
 - (ii) minimum average throughput (θ_{ave}, I), the lower bound of the throughput averaged over intervals of duration I ;
- (2) *delay bounds*: either
 - (i) deterministic maximum delay (D_{\max}): $D_i \leq D_{\max}$ for all i ; or
 - (ii) statistical maximum delay (D_{\max}, Z_{\min}): $Prob(D \leq D_{\max}) \geq Z_{\min}$.

The desired throughput bounds are specified implicitly in the client's characterization of the traffic to be carried by the real-time channel whose establishment the client is requesting. Calling "arrival time" the time a packet enters the network through the sender's interface, this characterization consists of the following parameters:

- (a) the minimum packet interarrival time x_{\min} ;
- (b) the minimum average packet interarrival time x_{ave} ;
- (c) the duration I of the time interval over which x_{ave} is computed; x_{ave} is the smallest of the averages, computed over all of the intervals of duration I , of packet interarrival times;
- (d) the maximum packet size s_{\max} ;
- (e) the service time t of a packet in each node along the channel's route.

Note that (d) and (e) are likely to be known to the network, and may not have to be specified by the client. Note also that requiring the client to specify x_{\min} and x_{ave} implies a pledge on the client's part to keep the minimum and average intervals between consecutive packets at the channel's entrance to the network larger than or equal to x_{\min} and x_{ave} respectively. The service provider must, however, ensure and verify that the client honors this pledge at all times, otherwise the performance given to other clients may be in danger. The solution to this problem adopted in [Ferr89] and [FeVe90] consists of having each node compute the minimum and average intervals, and suitably increase the local deadlines of those packets that do not satisfy the client's pledge. This *distributed rate control* algorithm is not only executed in the edge nodes, but in all nodes, since load fluctuations can indeed cause local and temporary violations even if all clients scrupulously honor their pledges.

The network we consider is a general-topology, point-to-point, store-and-forward, packet-switching one. The links are assumed to be bounded-delay: that is, the delay incurred by a packet over a link must have a finite and known upper bound. This requirement is satisfied by links consisting of pure communication channels, whose delays are those due to signal propagation, as well as by a number of more complicated structures, such as entire networks, but not by all of them; for instance, if the link connecting two nodes is an Ethernet, the bounded-delay requirement is not exactly satisfied.

Each node is assumed to be suitably modeled as a single server (in the queueing theoretical sense of the term "server") with multiple queues and multi-class deadline-based scheduling [Ferr89]. The node scheduler handles at least three queues, listed here in order of decreasing priority: the deterministic queue, where packets traveling on a real-time channel with deterministic delay guarantees are stored; the statistical queue, for packets on real-time channels with statistical guarantees; and a queue for everything else that is to be processed by the node. Each queue is ordered according to packet deadlines in the node. When the node is free to start shipping a new packet, the deadline of the packet at the head of the statistical queue is compared with the time at which the packet at the head of the deterministic queue should start being handled by the node for its local deadline to be met (we shall call this the "start deadline", to distinguish it from the "end deadline", or simply "deadline"). If the former is larger than the latter, the deterministic packet is shipped; if it is smaller, the statistical one takes precedence. Packets in the third queue are shipped only when the other two queues are empty. The algorithm can be generalized in the obvious way to the case of q queues, $q-1$ of which are for real-time packets with different priorities.

That of real-time channel is a logical concept: its physical realization could in principle be of various types, including multiple routes or datagram service between the sender and the receiver. However, obtaining performance guarantees and sequenced delivery (an important requirement in video and audio communications) is either very hard or impossible if a real-time channel is not realized as a connection with a single, fixed physical route. Since reservation of resources is necessary (not in the circuit-switching sense of this term, but to keep track of the potential real-time load on each resource and to prevent this load from exceeding its bounds), real-time channels are to be established explicitly when they are requested by a client (unless an existing one can be used for that client's traffic). The establishment procedure consists of selecting a route in a way analogous to that for a virtual circuit, though the routing algorithm might somehow be influenced by the knowledge that the channel to be set up is a guaranteed-performance one. However, besides selecting the route, the procedure will include computations and tests in each node along the route, the tentative reservation of node resources, and the transmission of information from each node to the receiver.

If the tests in a node are all successful, the *establishment message* is sent to the next node selected by the routing algorithm, and so on until it reaches the receiver. The receiver is the final decision-maker, i.e., it runs the last tests; if these too are successful, then the request is accepted, and the receiver sends an *acceptance message* back to the sender; this message transmits information from the receiver to each node on the route, and commits the resources (or the fractions of them that are really needed) tentatively reserved by the establishment message. If, on the other hand, one of the tests is unsuccessful, either in a node or in the receiver, the establishment message stops there, and a *rejection message* is sent back to the source. This message undoes the partial, tentatively established connection, and releases the tentatively reserved resources in each node along that connection.

The choice of the receiver as the final decision point was dictated by the desirability of a fast establishment procedure. Indeed, our procedure requires a single round-trip; it can, however, be made even faster by optimistically starting packet shipments before the new channel has been completely established [DaVe89].

If the request for the creation of a new real-time channel is accepted, each node along its route has to keep information about the channel; some of this information is to be used for the correct handling of the packets traveling on the channel, some other during the tests for the

establishment of new channels through the node, and some for both purposes. Besides storing routing information and all the channel's parameters specified by the client (e.g., x_{\min} , x_{ave} , s_{\max} , and so on), the node must be given and keep the values of the following parameters:

- (a) the *delay bound* $d_{i,n}$ of channel i in node n ; this is the maximum (in a deterministic or statistical sense, depending on channel i 's type) time a packet on channel i should spend in the node; the *deadline* of a packet arriving at the node at time t will be computed as $t+d_{i,n}$, unless the rate control algorithm (see Section 3) intervenes to modify it;
- (b) the *delay bound compliance probability* $z_{i,n}$ of statistical channel i in node n ; this is the minimum probability that a packet on channel i will not be delayed more than $d_{i,n}$ in node n .

The values of $d_{i,n}$ and $z_{i,n}$ are computed by the receiver, the first entity, among those involved in a channel's establishment, that knows all the necessary information to do these computations. While $d_{i,n}$ is used both to determine packet deadlines and in the delay bound test for the establishment of a new channel (see below), $z_{i,n}$ only appears in establishment tests (see the statistical test below). Since the network has an arbitrary topology, each node must have its own independent and local characterization of the guarantees it has given each channel passing through it. This role is played by $d_{i,n}$ and $z_{i,n}$, and is the only justification for the latter parameter's existence. For the same reason, the value of $Z_{min,i}$ specified by the client is factored by the receiver into the individual node compliance probabilities $z_{i,n}$:

$$Z_{min,i} = \prod_{n=1}^N z_{i,n} , \quad (1)$$

where N is the total number of nodes along the channel's route. This implies the pessimistic assumption that a packet delayed beyond its local bound in a node will not be able to satisfy the channel's global delay bound.

The tests to be performed in each node are as follows (see [FeVe90]).

(1) *The deterministic test.* This test is to be performed for any deterministic channel to be established, and involves the deterministic channels already passing through the node. It verifies that the node has enough processing power to handle the additional deterministic channel without impairing the guarantees given to the other channels in the node; the power is sufficient if the utilization of the node's processor due to deterministic channels is less than 1 even after the addition of the new channel and even in the worst case. The computation requires knowledge of t and x_{\min} for each of the channels involved. Note that we can assume that x_{\min} and x_{ave} remain the interpacket time bounds along the entire route of the channel because of our rate control policy (see Section 3).

(2) *The statistical test.* This is to be performed whenever at least one statistical channel is involved, i.e., either to be established or already established in the node. It determines whether, for each statistical channel j already passing through node n , the probability P_{do} that its packets will incur a delay higher than its bound $d_{j,n}$ in the node is below the maximum tolerable value for that probability, $1-z_{j,n}$. P_{do} is assumed to be the same for all packets on all channels passing through the node, and is computed from the probabilities of all the combinations of active channels that can potentially saturate the node. This computation assumes that channel traffic patterns are statistically independent. The value of

P_{do} thus calculated is transmitted to the receiver via the establishment message, to be used in the computation of $z_{i,n}$.

(3) *The delay bound test.* This test is performed in all cases, to determine whether *scheduler saturation* can be avoided in the node even after the creation of the new channel. Scheduler saturation makes it impossible for the node scheduler to satisfy all delay bounds in spite of the absence of *node saturation*, which is handled by the statistical test. This condition is caused by delay bounds that are too small compared with the service times t and the minimum interpacket times x_{\min} . To perform the delay bound test, we compute the minimum delay bound $d_{i,n}^l$ that can be assigned to the new channel in the node. This computation requires the knowledge of $x_{\min,j}$, $t_{j,n}$, $d_{j,n}$, $x_{\min,i}$, and $t_{i,n}$ ($j \neq i$). If the new channel passes the test, the value of $d_{i,n}^l$ is transmitted to the receiver via the establishment message.

The receiver performs two tests and two computations to determine the values of $d_{i,n}$ and $z_{i,n}$, which it then sends to all the nodes on the route via the acceptance message. If the channel to be established is deterministic, only one test and one computation (that of the $d_{i,n}$'s) are needed. The tests are called the D test and the Z test.

(1) *The D test.* This test verifies that the sum of the $d_{i,n}^l$'s over the channel's route is smaller than or equal to $D_{\max}^* = D_{\max} - \sum \text{link delays}$.

(2) *The Z test.* This test consists of verifying that the combination of all delay bound compliance probabilities will be greater than or equal to $Z_{\min,i}$:

$$\prod_{n=1}^N (1 - P_{do,n}) \geq Z_{\min,i}. \quad (2)$$

There are various possible approaches to the computation of $d_{i,n}$ and $z_{i,n}$. One of them, for example, is described in [FeVe90].

3. A Buffer Space Allocation Algorithm

In [FeVe90], another client requirement the real-time channel establishment scheme could satisfy, besides the throughput and delay bounds defined in Section 2, was a *reliability bound* (W_{\min}):

$$\text{Prob}(\text{packet sent is delivered correctly}) \geq W_{\min}. \quad (3)$$

There are two possible causes of packet loss:

- (i) transmission error, usually due to noise or failure of some network component, and
- (ii) buffer overrun, usually due to a space shortage in a node of the network or in the receiver.

While problem (ii) causes packets to be really lost, i.e., to go undelivered, transmission errors often garble them. In real-time wide-area communication, the use of such error control mechanisms as acknowledgements and retransmission is highly questionable, since the long delays inherent in the operation of these mechanisms are likely to prevent the packets eventually delivered correctly from satisfying their delay bounds, hence from being usable by the receiver. Thus, we can, at least in most applications, consider as lost any packet that is delivered incorrectly.

Of all the packets that are shipped on a real-time channel, at least $100W_{\min}\%$ must be delivered correctly. If the channel is statistical, at least $100W_{\min}\%$ of those delivered correctly must be delivered on time. Thus, the packets that are delivered correctly and on time will be at least $100W_{\min}Z_{\min}\%$ of those that were shipped by the sender. Some clients may be more interested in $W_{\min}Z_{\min}$ than in the individual values of these two bounds, and may want to leave the decision about what these values should be to the real-time server [Ferr90].

We assume that the network is characterized by a known error rate, and that the given value of W_{\min} will be divided by the complement to 1 of this error rate to obtain the lower bound for the probability of buffer overrun. To avoid complicating our symbology, we shall use W_{\min} in the sequel to denote not just the correct-delivery probability bound, but one of its two factors, namely, the lower bound of the probability that no buffer will overflow, thereby causing packet loss.

In the scheme presented in [FeVe90], the probability of packet loss was always 0, since we assumed no limitations on the amount of buffer space allocated to each channel in each node. In this paper, we remove that assumption, and propose and evaluate an algorithm that assigns reasonable amounts of buffer space to channels taking into account the given value of W_{\min} . This value will be factored by the receiver into node contributions $w_{i,n}$ as is done for Z_{\min} and because of the same reasons explained in Section 2:

$$W_{min,i} = \prod_{n=1}^N w_{i,n}. \quad (4)$$

The implication of this expression, that is, that a packet lost in node n is lost forever, is of course true. The treatment of the problem for statistical channels is quite different from (and much more complex than) that for deterministic channels. After examining both, we briefly describe an algorithm, based on these treatments, for computing the space to be allocated in each node to a new channel being established.

3.1. Buffer space for a deterministic channel

Let us assume $w_{i,n} = 1$ for channel i in node n . This means that the client has requested a guarantee that there will be no packets lost because of buffer overruns on that channel. The scheme described in [FeVe90] guarantees that in node n , if channel i is deterministic, each of its packets will spend a time not longer than $d_{i,n}$. In the worst case, a packet on channel i will arrive at the node at intervals of length $x_{min,i}$. Assuming the residence time of each packet in node n to be $d_{i,n}$, which is clearly a worst-case assumption, the amount of space in bytes needed by the channel in that node will be

$$SPACE(i,n) = s_{max,i} \left\lceil \frac{d_{i,n}}{x_{min,i}} \right\rceil. \quad (5)$$

If we now assume $w_{i,n} < 1$, we may be able to reduce $SPACE(i,n)$ by up to u packet buffers. In the worst case, subtracting u buffers from the $b = SPACE(i,n)/s_{max,i}$ buffers that are needed to avoid losses may cause up to u packets out of every b to be dropped. Thus, we will have

$$1 - w_{i,n} = \frac{u}{b}, \quad (6)$$

hence the number of packet buffers needed will be

$$b' = \lceil b - u \rceil = \left\lceil bw_{i,n} \right\rceil, \quad (7)$$

and the amount of buffer space in bytes

$$SPACE(i,n) = s_{max,i} \left\lceil \frac{d_{i,n}}{x_{min,i}} w_{i,n} \right\rceil. \quad (8)$$

3.2. Buffer space for a statistical channel

We begin by assuming again $w_{i,n} = 1$. The scheme described in [FeVe90] is such that, if node n is not saturated (i.e., if $\sum(t_{i,n}/x_{min,i}) = 1$), all packets carried by any statistical channel satisfy that channel's delay bound $d_{i,n}$. Thus, a non-saturated node treats deterministic and statistical channels in the same way (though the actual delays incurred by the former are shorter because of the scheduler's preference for them). When, on the other hand, the node is saturated, some of the statistical packets may be delayed beyond their bound and have therefore to stay longer in the buffers. In saturation, expressions (5) and (8) no longer apply.

Let $d_{i,n}^{sat}$ be the maximum delay bound in node n for channel i in a particular saturated situation. It should be noted that to each such situation there will generally correspond a different value of this bound; however, given certain channels established in node n , with $\sum(t_{i,n}/s_{min,i}) > 1$, there will be a maximum packet residence time in the node, and any value greater than or equal to this time will be a valid one for $d_{i,n}^{sat}$.

The maximum time spent by a channel i packet in node n can be computed by observing that the worst case will occur at the end of the longest possible interval during which all channels are simultaneously active and are sending packets at their fastest rates. Let t_a be an instant at the beginning of this interval when a packet on channel i arrives at node n . Since a relatively long time usually elapses between the instant a node becomes saturated and the first instant at which a delay bound is violated, we may assume that this packet leaves the node at a time not later than $t_a + d_{i,n}$. Let now t_b be the arrival time of the channel i packet closest to the end of the saturation interval being considered; this packet will leave the node at time $t_b + d_{i,n}^{sat}$, and the node will be busy without interruptions at least between $t_a + d_{i,n}$ and $t_b + d_{i,n}^{sat}$.

The packets that cause this activity in the node are those whose (start or end) deadlines are greater than $t_a + d_{i,n}$ and smaller than $t_b + d_{i,n}$. During an interval of duration $t_b - t_a$, the number of these packets will equal the total number of arrivals, and each arrival on a channel j will contribute a term $t_{j,n}$ to the node's busy time. Thus, we can write

$$t_b + d_{i,n}^{sat} - t_a - d_{i,n} = \sum_j \left\lceil \frac{t_b - t_a}{x_{min,j}} \right\rceil t_{j,n}, \quad (9)$$

from which, ignoring the ceiling and noting that the portion of the worst interval of duration I in which node n is saturated has length $I P_{do}$, we obtain

$$d_{i,n}^{sat} = d_{i,n} + I P_{do} \left(\sum_j \frac{t_{j,n}}{x_{min,j}} - 1 \right). \quad (10)$$

To calculate $SPACE(i,n)$ when n is saturated (i.e., when $\sum(t_{j,n}/x_{min,j}) > 1$), we replace $d_{i,n}$

in (5) with $d_{i,n}^{sat}$ as given by (10). When n is not saturated, expression (5) will suffice even for a statistical channel.

If $w_{i,n} < 1$, we can use (8) in a non-saturated node, and

$$SPACE(i,n) = s_{max,i} \left[\frac{d_{i,n}^{sat}}{x_{min,i}} w_{i,n} \right] \quad (11)$$

in a saturated node.

3.3. The algorithm

If node n is not saturated, the amount of buffer space a real-time channel i to be established needs in n can be determined using equations (5) or (8), depending on whether $w_{i,n} = 1$ or $w_{i,n} < 1$. If i has a statistical delay bound and n is saturated, then we must use either equation (5) with $d_{i,n}^{sat}$ (given by (10)) in lieu of $d_{i,n}$ or equation (11), again depending on whether $w_{i,n} = 1$ or $w_{i,n} < 1$.

In all cases, we need to know the value of $d_{i,n}$. Since during the forward trip of the establishment message this value is not known, we have to estimate it in order to determine whether the node has sufficient buffer space for the new channel, and, if so, to reserve that space until an acceptance or rejection message concerning channel i is received by the node. Obviously, since we cannot compute $d_{i,n}$ exactly, we have to overestimate it. A simple, though very gross, estimate satisfying this condition is the value of the channel's delay bound $D_{max,i}$; of course, better (but more expensive) estimates can also be easily devised. An acceptance message carries within itself the value of $d_{i,n}$, which is used by the node to compute the amount of buffer space to be permanently reserved for channel i ; if it was grossly overestimated, this amount will be substantially smaller than the one that was set aside during the forward trip.

If $W_{min,i} < 1$, it is likely that we will have $w_{i,n} < 1$. Since the exact value of $w_{i,n}$ is not known during the forward trip, we will set $w_{i,n} = 1$ in (8) and (11) at that time; in other words, we will always use (5) or its version with $d_{i,n}^{sat}$ during the forward trip, and switch to either (8) or (11) if and when an acceptance message is received by the node.

If node n is saturated and channel i is statistical, we use (10) to compute $d_{i,n}^{i,n^{sat}}$. The term $\sum t_{j,n}/x_{min,j}$, which must be computed in any case (even if the deterministic test is not needed) to determine whether n is in saturation, includes the contribution of channel i . The value of P_{do} is known from the statistical test computation, which must therefore be performed before the *buffer space test*. Again, we do not know $d_{i,n}$, and we have to estimate it, for example by setting it equal to $D_{max,i}$ during the forward trip.

Since every new channel established through the node increases the value of $d_{i,n}^{sat}$ for each of the channels already traversing the node, $SPACE(j,n)$ must be computed at every new request for each channel j when and until the node is saturated, and will normally grow. If the space allocated to channel j was equal to the value of $SPACE(j,n)$ calculated when the node was not saturated, then either we are allowed to increase space allocations or we have to reject the request. Thus, if space increases for existing channels are to be excluded, the initial allocation should be greater than $SPACE(j,n)$, otherwise the potential for a better exploitation of the network's resources by statistical channels will be severely curtailed.

If channel space allocations are allowed to grow, it would seem useful to reduce them when existing channels in the node are deleted. However, this is usually not a wise strategy, as the space freed by the deleted channels is likely to suffice for the subsequent creation of new channels, whose number is going to be limited by the same admission control mechanisms that have been operating since the beginning. In other terms, the buffer space a channel has obtained up to that point is a better estimate of its needs at that level of node saturation than any of the previous allocations. The only persuasive reason for reducing the sizes of space allocations would be to give some more space to non-real-time traffic if this has been squeezed too much by the growth of the real-time component. Note that this would be a symptom of a design flaw, i.e., of insufficient total buffer space in the node.

After the amount by which the allocations should grow has been determined, the buffer space test simply consists of checking whether there is enough free space for this growth. Within limits, this "free" space may come from the pool that is being used by non-real-time traffic.

The receiver, together with the values of $d_{i,n}$ and $z_{i,n}$, computes those of the $w_{i,n}$'s, for example as follows:

$$w_{i,n} = W_{min,i}^{1/N}, \quad (12)$$

where N is the total number of nodes along the channel's route. A better approach is one that does not reject a request in a node n if the space available is smaller than the estimated value of $SPACE(i,n)$, but introduces the notion of *offered buffer space*, defined as $B_{o,n} = \min(\text{available space}, SPACE(i,n))$. If we set

$$b_{i,n} = \frac{B_{o,n}}{SPACE(i,n)}, \quad (13)$$

and send $B_{o,n}$ and $b_{i,n}$ to the receiver, the receiver may compute $w_{i,n}$ as

$$w_{i,n} = \left(\frac{W_{min,i}}{\prod_n b_{i,n}} \right)^{1/N} b_{i,n}, \quad (14)$$

which assigns a smaller value of $w_{i,n}$ to those nodes where the offered buffer space is smaller than the required one. Of course, channel i is rejected if (14) yields $w_{i,n} \leq W_{min,i}$ for some n . This may be called the *W test*, to be performed by the receiver along with the D test and the Z test (see Section 2).

It should be noted that the algorithm described in this section and the equations on which it is based assume that the minimum interval between the arrivals of two consecutive packets of a channel i at a node is not smaller than $x_{min,i}$. Unfortunately, the distributed rate control scheme referred to in Section 2 does not always guarantee that this assumption will always be satisfied. To solve this problem, we can resort to one of the following approaches.

(1) The maximum number of packets a node, suddenly finding itself with only packets of one channel, can send to a neighbor in rapid succession is bounded by the space allotted to that channel in that node. Hence, if we assign to channel j in node n an amount of space equal to the value of $SPACE(j,n)$ computed as illustrated above plus the space allotted to j in the previous node, there can never be an overflow. However, this approach has the undesirable properties of being quite wasteful and of requiring amounts of space that are

rapidly growing from source to destination.

(2) The node scheduling algorithm could be modified so that no packet will be shipped by a node before the time it would have arrived if it had not violated the maximum rate bound. This would allow the additional space to be given to channel j in node n to be limited to $SPACE(j, n-1)$ instead of $\sum_{k=1}^{n-1} SPACE(j, k)$. Thus, at the expense of a slightly more complicated scheduler, the waste would be reduced, and the space allotments along a path would be stable rather than growing.

(3) Instead of being purely static, the allocation of space could be partially static and partially dynamic. In other terms, besides allotting $SPACE(j, n)$ bytes to channel j in node n , we could maintain a small dynamically managed buffer pool to absorb the effects of upstream load fluctuations. The maximum amount of dynamic space a channel could occupy would have to be limited to prevent the pool from being monopolized by faulty or malicious clients. This bound could be a function of the amount that would have to be statically allocated to the channel with approach (1) or (2) above, depending on the node scheduling algorithm.

4. The Simulations

In this section, we attempt to provide simulation-based answers to the following questions which the reader may have asked about the scheme.

- Is the scheme correct? Can we discover cases when the loss guarantees cannot be met using this scheme?
- How much time is required to perform the buffer tests?
- Providing the loss rate guarantees involves an allocation of buffers to each channel. What are the actual buffer utilizations for these channels? How do they compare to the ones predicted by equation (11) and equation (8) ?
- Can we identify the traffic conditions and performance requirements under which the scheme performs well, and, conversely, the circumstances under which it fails to perform well?
- How do the deterministic and statistical channels compare relative to each other? What is the space cost of a statistical channel as compared to the space cost of the deterministic channel with similar traffic characteristics?

We obtained answers to these questions by means of a simulator written in CSIM [Schw89]. In this section, we describe simulations of two networks, one a very simple network consisting of two nodes only, and a complex network of 14 nodes with a mix of different types of traffic channels to illustrate the feasibility of the algorithms for different traffic types and arbitrary topologies. The workload components chosen were the same as for the simulations of the first part of the establishment scheme, which dealt with delay bounds only [FeVe90]. For the reader's convenience, we describe the network and the workload components again in this paper.

Figure 1(a) shows the graph that was used in the simulations. Nodes 0-8 serve as sources and destinations for the real-time channels. These nine nodes are divided into groups of three nodes each. Each source generated channel requests at random intervals. The destination group of each channel request was chosen by the probability matrix in Table 1. The actual destination

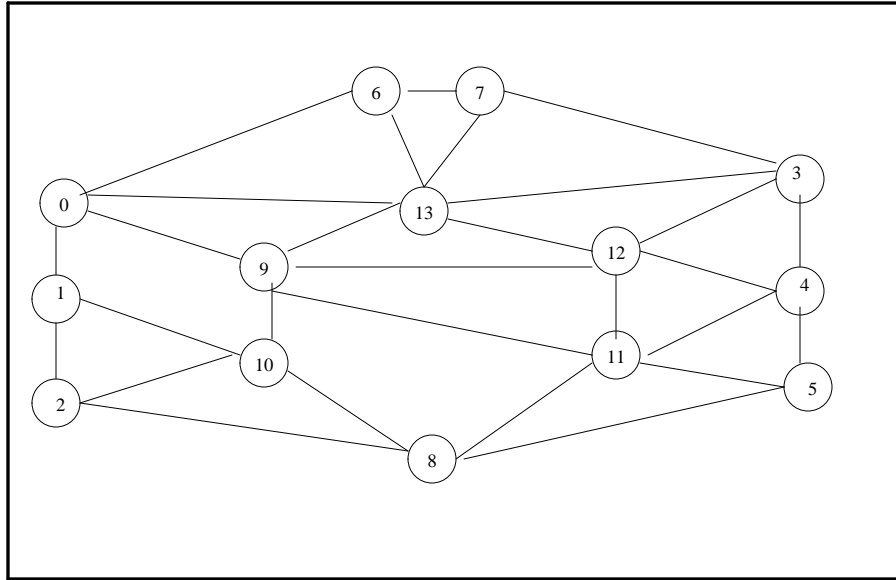


Figure 1(a). The network that was simulated. Nodes 0-8 serve as sources and destinations. Nodes 9-13 serve as gateways.

was picked randomly from the group. One of the four shortest paths from a channel’s source to its destination was picked at random as the route of its establishment request.

Table 1. Source destination group probabilities			
From nodes	To nodes		
	0-2	3-5	6-8
0-2	0.10	0.72	0.18
3-5	0.72	0.10	0.18
6-8	0.45	0.39	0.16

4.1. Workload

Future multimedia networks will have a wide variety of traffic types. To capture some of this diversity, we consider four types of channels. Table 2 shows the types of channels used in the simulation experiments. Type I channels tax the network resources heavily, while type IV channels are very light.

Table 2. Channel types			
Channel Type	t (time units)	x_{\min} (time units)	x_{ave} (time units)
I	4	10	60
II	4	40	120
III	1	10	60
IV	1	40	120

The simulation was run in two phases. In Phase 1, we tried to establish as many channels as we could. The next channel to be established was of any type and either deterministic or statistical with equal probability. The exact number and types of channels set up in such a context depended on the order of arrival of channel requests. In some cases, a request for a type I channel arrived first, and many subsequent requests for type IV channels were rejected. On the other hand, if the first few channels were of type IV, a subsequent request for a type I channel might be refused, but many more channels of type III could be established.

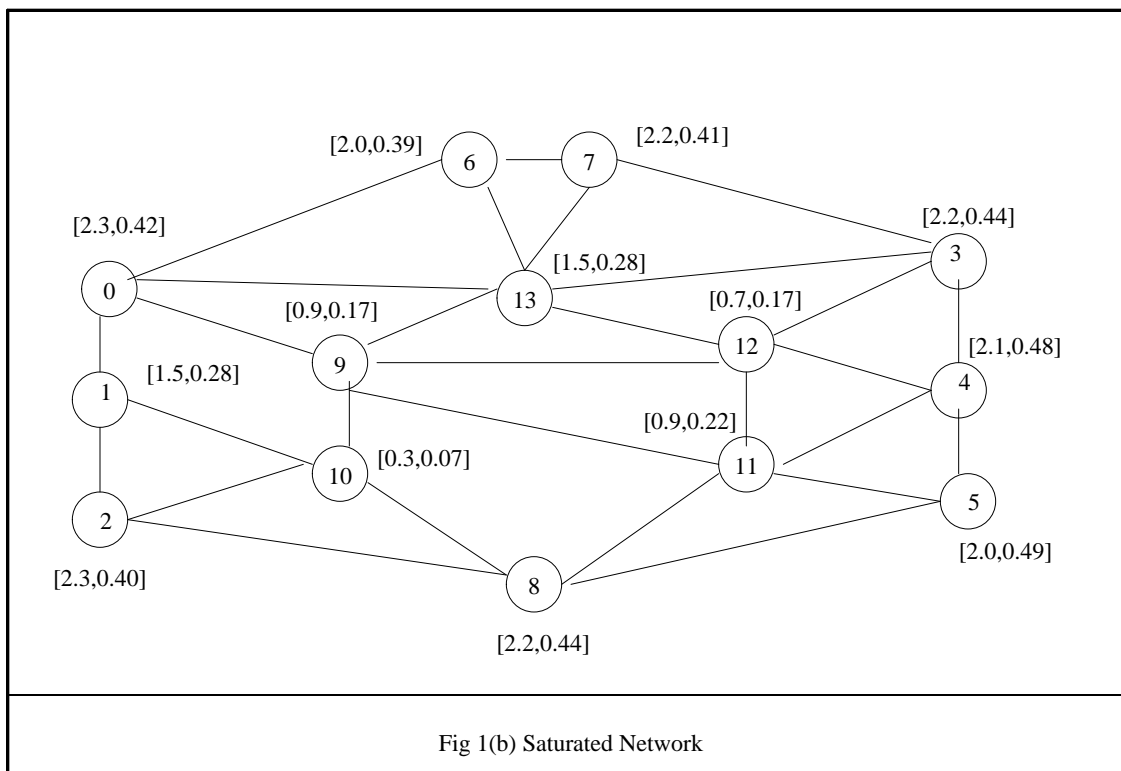


Figure 1(b). Result of a Phase 1 simulation of the network in Figure 1(a). The numbers reported near each node are the sum of the weights of all established channels at the node and the expected utilization of the node by these channels. The numbers near node 3 are the most interesting since the delay and buffer distribution at node 3 have been plotted in Figure 1(c).

Let the *weight* of a channel be defined as t/x_{\min} . A correct measure of saturation at a node is the total weight, rather than the actual number of channels accepted.

The results of one Phase 1 simulation of the network in Figure 1(a) are shown in Figure 1(b). The two numbers associated with each node are the weights of all the channels and the expected utilization of the node by the established channels. We selected this configuration, containing about 56 successfully established channels for Phase 2 simulations. Note that routing constraints prevent some nodes, e.g. node 12, from being saturated in this network¹.

¹ The set of channels accepted in Figure 1(b) is somewhat different from the configuration shown in [FeVe90]. However, the peak weights and utilizations of the nodes are almost identical. Both these graphs show a typical case of node utilization. The main reason for these slight differences are that we modified the statistical test to use a much faster, although slightly pessimistic approximation for evaluating the overflow probability. Reproducing the

In Phase 2, we experimented with a number of different arrival patterns. The worst pattern, i.e., the one which resulted in the maximum delays and losses, appeared to be a pattern in which packets arrived either at an interval of x_{\min} or at an interval x_l that brought the average in a period of duration I to x_{ave} . Thus, the arrival pattern consisted of a burst of $\lceil I/x_{ave} \rceil$ packets followed by a period of silence. In order to satisfy the independence assumptions involved in the computation of the overflow probability, all the traffic patterns were given a random offset in every interval of length I , which staggered them relative to one another.

We also experimented with other arrival patterns, for example, one in which the burst sizes were distributed geometrically, one in which arrivals occurred smoothly at the average rate and one in which the inter-arrival time was uniformly distributed between the minimum and an upper limit so as to achieve the specified average value. However, we will only describe the results of the worst arrival pattern, i.e., one in which the arrivals consisted of the maximum possible burst size followed by a long period of silence.

4.2. Results

In all the various experiments we ran, we were unable to find any situations where the delay or loss guarantees were not satisfied. Figure 1(c) shows the delay distribution of a statistical channel at node 3 (which had a utilization of 44%). The channel was of type III and had a delay bound of 39 units at node 3 with z (probability of meeting this delay) equal to 0.90 and w (probability of not being lost) equal to 0.95. The value of I for the channel was 920 units, while the value of P_{do} , the node overflow probability was 0.05, due to presence of other statistical channels. The formulas in Section 3.2 tell us that 4 buffers would have caused a loss rate of no greater than 5% since the probability of overflow was about 5%. In the simulation, however, 7 buffers were allocated because of our decision to overallocate on the forward pass of the channel establishment request, and the decision not to reduce the buffer space allocation even in the case of channel deletion.

Figure 1(c) shows that a packet is likely to have a very low delay although for some packets the delay may be as high as twice its assigned bound. In the simulations, 98.8% of the packets were transmitted within the required bound of 39 units. Figure 1(d) shows a histogram of the number of packets present in the queue (or buffer occupancy) as seen by packets arriving at node 3 on the same channel to which Figure 1(c) refers. Most arrivals find no packet from the same channel in the queue. We had allocated 7 buffers for this channel, and in the actual simulations no packets were lost due to buffer overruns. Notice that about 93% of the packets found the buffer space to be empty in Figure 1(d), but only about 72% in the node had a delay equal to the service time of one packet. There is no contradiction in these figures, however since there may be packets present from other channels which cause the channels under consideration to have slightly higher delays.

Thus, Figures 1(c) and 1(d) show that the aforementioned formulas and the algorithm are sufficient to meet the required performance bounds for a statistical channel.

In the case of a deterministic channel passing through the same node, our simulations did not show any loss of packets. Thus, both deterministic and statistical guarantees were met by our scheme.

exact set of channels was difficult, but the typical node utilizations and delay distributions remained similar.

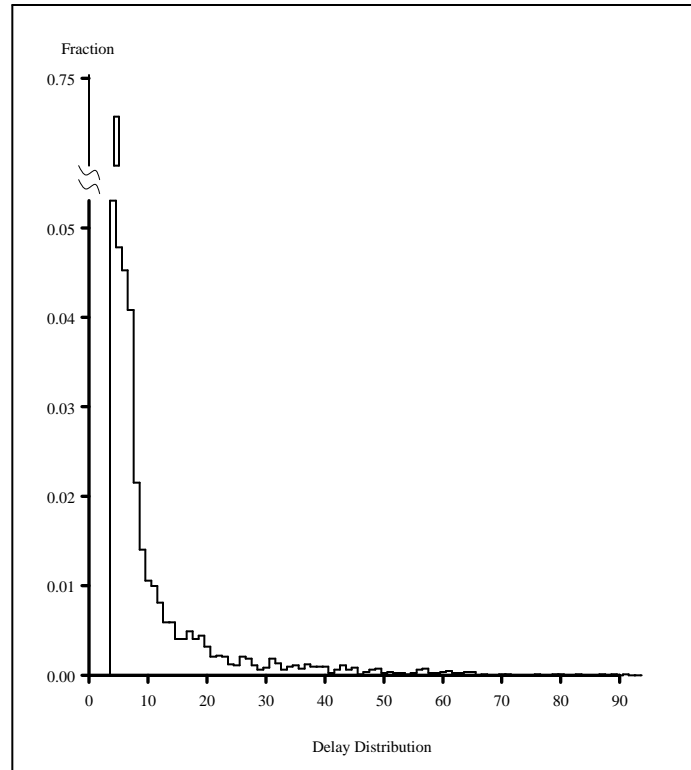


Figure 1(c). The delay distribution on a statistical channel. A large fraction of the packets had a very low delay (which is why there is a break at the top left corner of the histogram). The delay bound of the channel was 39 units, and about 1.4% of the packets failed to meet the required delay bound. The maximum delay, however, could be substantially higher than the assigned bound, e.g. in this figure the maximum delay is 78 units, about twice the bound of 39 units.

The overhead of the buffer test is linear in the number of channels already established at a node. Similarly, the overhead of the destination algorithm is linear in the number of the nodes along the path. To give a quantitative estimate, with between 8 and 18 channels established at various nodes, the mean time to perform the buffer tests was about 1 milliseconds at every node. The run time overhead for buffers consisted of a simple increment and comparison at packet arrival and a simple decrement at packet departure, which did not take any significant amount of time.

4.3. Simulations of the simple network

Having convinced ourselves that the scheme does work, let us try to answer questions about the relative cost of statistical and deterministic channels, and about the scheme's performance. In order to study the effect of various parameters like W , Z and I on the amount of buffer space required, we chose a very simple network, one consisting of only two nodes and studied the different aspects of our establishment scheme using channels of only one type. We arbitrarily chose channels of Type III for these studies.

Figure 2(a) shows the amount of buffer space required by our scheme for a statistical channel. Three sets of curves are shown, for value of I corresponding to 480, 960 and 1920 units,

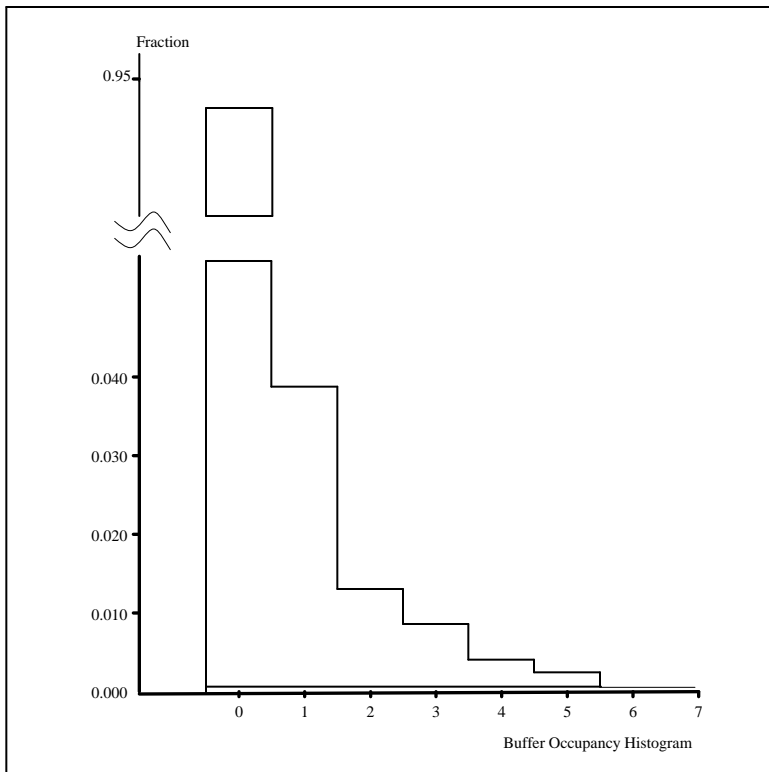


Figure 1(d). The buffer occupancy distribution for a statistical channel. The figure shows the histogram of the number of packets of the same channel present in the channel's queue at node 3 as seen by arriving packets. 7 buffers were allocated to this channel. This bound appears to be sufficient since all 7 buffers were used at some time in our simulation run. However, no packets were dropped even though the space allocated allowed us to drop 4% of the packets. If we had allocated only 4 buffers, then also we could have met the required performance guarantees, as predicted by the formulas in Section 3.2. A large fraction of the packets find no buffers occupied, as shown by the break in the histogram at the top left corner.

respectively. Within each set, the curves correspond to a local delay bound d of 11, 22 and 33 units respectively. Recalling that channels of type III have a x_{\min} of 10 units, a value of d less than 10 would have resulted in the node rejecting channels due to the delay bound test (see Section 2) before node saturation can occur. Thus, the overflow probability would have been zero and buffer allocation would be done by equation (5) in Section 3.1.

For the case of higher delays, where the statistical test is the cause of channel rejection (and thus the bottleneck), we see that buffer allocation is much more sensitive to the averaging interval I than it is to the local delay bounds assigned. Notice also that the curves reach an asymptote which seems to depend only on the value of I . The explanation for the phenomenon is simple. As the value of overflow probability increases, the buffer space predicted for statistical channels by equation (11) increases. However, the traffic constraints of the channel themselves impose an upper limit on the total number of packets that can ever be present at a node. Since an interval I can have at most $\lceil I/x_{\text{ave}} \rceil$ packets, this is a natural upper bound on the number of packets that can ever be present at a node. Thus the number of buffers required for very high values of the overflow probability is constant².

² This is true only if we assume that the node utilization is less than unity. On the other hand, this condition is a basic requirement in order to provide any kind of performance guarantee, since neither losses nor delays are bound-

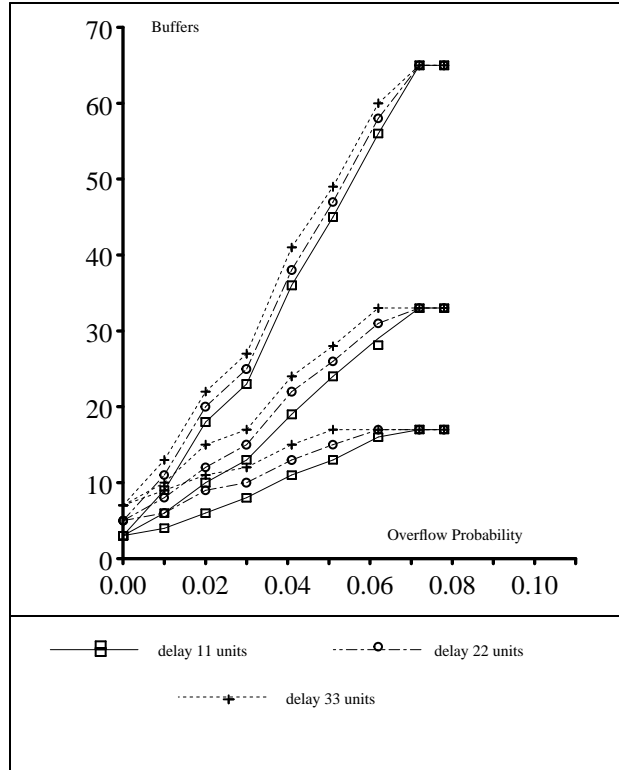


Figure 2(a). The effect of the overflow probability on the amount of buffer space required for statistical channels. The three sets of curves are for type III channels (see Table 2) with different values of the averaging interval I . Within each set, the three curves correspond to a local delay bound of 11, 22 and 33 time units respectively. The values have been chosen to avoid scheduler saturation (see Section 2).

One result which is obvious from Figure 2(a) is that a statistical channel may have substantial storage costs associated with it as compared to a deterministic channel. This is especially true for bursty traffic, which may require a higher averaging interval than a smooth one.

Figure 2(b) shows the effect of yet another parameter, the loss probability $1 - w_{i,n}$ at a node. The delay bounds of all the channels have been chosen to be 11 time units in every node. This is the minimum required to avoid scheduler saturation. The value of the overflow probability at each node is 0.08. Notice that in the case where scheduler saturation occurs, the overflow probability at a node is zero and hence equation (5) would become the buffer allocation formula; the result of this formula in this case would be 2 buffers. One interesting aspect of the graph is the convergence of the curves to the deterministic allocation at high values of loss probability. This convergence occurs whenever the value of the loss probability exceeds the overflow probability. In Figure 2(b), it occurs at values roughly greater than 0.08.

The result is intuitively sound since in the case where the loss requirements are less stringent than the delay bound requirements, it makes sense to discard any packet that might get delayed. Our buffer allocation policy achieves exactly the same effect. In the case of node saturation, when packets are likely to be delayed, arriving packets on such a channel are

ed in a node with an average utilization exceeding unity.

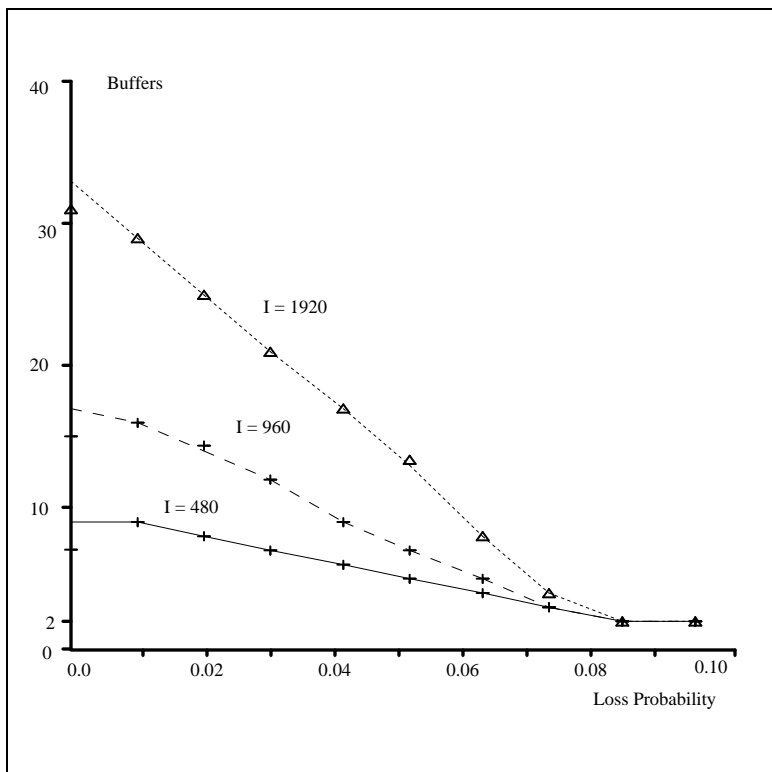


Figure 2(b). The effect of loss probability on the amount of buffer space required for statistical channels. The three curves are for type III channels (see Table 2) with different values of the averaging interval I and with a local delay bound of 11 time units, just sufficient to avoid scheduler saturation. Notice the convergence to the deterministic allocation (equation (5)) for higher values of loss probabilities.

discarded. A smarter gateway may accept the new packets and overwrite the existing packets which would have a higher delay. In both these cases, the performance guarantees of the channel are being met anyway.

Let us now examine the utilization of the buffer space allocated to the channel. To some extent, a partial answer is provided by Figure 1(d) which shows that all the seven buffers allocated for the channel shown there were used. However, we can also notice that the channel client there allowed the network to lose as much as 5% of the packets, while none were lost; and if we had allocated only 4 buffers, i.e. given up on the policy of not reducing buffers already allocated to a channel, they would have been sufficient to keep the loss rate lower than 2%. Under what conditions will this good correspondence be observed? To answer this question, we looked at the actual utilization of buffers in the two node network with channels of type III, and gave up the policy of not reducing the overallocated buffers to an established channel. Figure 2(c) shows that the actual utilization of buffers allocated to a deterministic channel is very high. In these simulations, one buffer was not used, but that was more an artifact of using the ceiling of $11/10$ (d/x_{\min}), which allocated two buffers while only one was used, and an extra buffer was allocated to cushion the effects of distributed rate control. The waste therefore is not substantial, at most one buffer, and if the delays are significantly higher than the value of x_{\min} , the utilization is quite high.

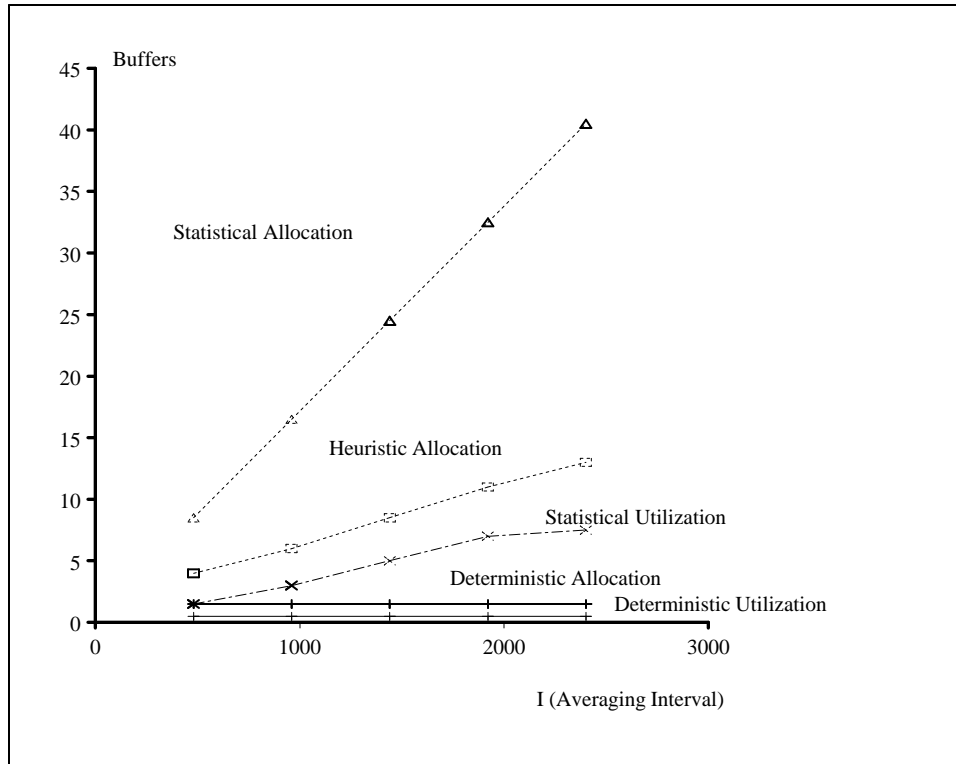


Figure 2(c). The actual utilization of the buffers allocated. The two allocation curves are for a deterministic channel of type III with a local delay bound of 11 units. The statistical channel has the same delay bound but an overflow probability of missing the deadline as 0.08 at each node. The allocation is done according to equation (5) for deterministic channels and according to equation (11) for the statistical channels. The utilization curves show the actual number of buffers ever occupied by a channel in phase 2 simulations. The allocation of statistical channel does seem to be too conservative. However, some heuristics may be used to reduce the allocation.

The actual utilization for the statistical case is not too good. Only a small fraction of the buffers allocated is used and the situation worsens as the averaging interval I increases. However, it is not clear whether a formula tighter than equation (11) can be derived using only the information used in Section 3. However, one can try using some heuristics to decrease the actual allocation to reasonable limits. One such heuristic is to use a smaller value of w hoping that there will be no overflows. Such a heuristic, with w equal to 0.92 is also shown in Figure 2(c). Again, the heuristic does not seem to work too well for large values of I . Moreover, there is no certainty that it will be sufficient for small values of I . Perhaps our simulation runs did not discover the worst possible scenarios, and then any guarantees verified by such heuristics would be suspect. Another approach would be to give up a purely static allocation and use a dynamic pool of buffers for statistical channels, in a manner similar to the buffer space required for distributed rate control (see Section 3.3).

It would thus appear from Figures 2(a), 2(b) and 2(c) that the scheme performs best for smaller values of the averaging interval I , and is sensitive to the delay bound compliance probability (to a limited extent) and to the reliability bound. Buffer utilization in the case of deterministic channels and when the reliability bounds are close to the delay bound compliance probability is reasonable. However, very large values of I tend to cause over-allocation of buffers.

Some heuristic needs to be used in such cases.

5. Conclusions

This paper has presented an algorithm for buffer space allocation in the establishment of real-time channels in a packet-switching network. The algorithm was the only part missing in the scheme for real-time channel establishment described in a preceding paper [FeVe90]. With the addition of the mechanisms illustrated in this paper, the scheme can therefore be regarded as complete, and serve as the basis for the design of a real-time communication service in a packet-switching network. The validity of our conjecture that it is feasible to build such a service on top of pure packet switching is therefore confirmed.

The buffer space allocation algorithm has been subjected to extensive simulation-based testing. The conclusions one can draw from those tests are that:

- the algorithm seems to be correct; i.e. our simulations have failed to reveal a case in which the packet loss guarantees could not be met;
- the overhead of the buffer-related computations during channel establishment is very reasonable;
- the utilizations of buffers allocated to deterministic channels are good, whereas those of statistical channels are often less good, and should perhaps be improved by replacing our worst-case formulas with suitable heuristics;
- the algorithm results in higher space utilizations when the channels are less bursty;
- since the delay bound of a deterministic channel is absolute, while some of the packets on a statistical channel may be substantially delayed, the buffer space needed for the former is smaller (sometimes much smaller) than the one for the latter, for equal delay and packet loss traffic characteristics.

The recognition that the algorithm discussed in this paper completes the work in [FeVe90] does not imply that no work remains to be done on this topic. In the area of the algorithm itself, a reliable heuristic that will reduce the space allocated to statistical channels without endangering the loss probability guarantees ought to be devised, and the impact on the mechanism of a non-deadline-based scheduling policy ought to be determined. Also a better understanding of the relationships among the various client requirements (in particular, those between statistical delay bounds and packet loss probability bounds), and their influences on the cost of real-time channel communications, should be reached for the benefit not only of the service designers but also of the clients[Ferr90]. Finally, the buffer space algorithm should be incorporated into the detailed design of a real-time service for packet switching networks.

6. References

- [CoYa88] D. E. Comer and R. Yavatkar, "Flows: Performance Guarantees in Best Effort Delivery Systems", *Proc. 8th Int. Conf. Distrib. Comp. Sys.*, San Jose, CA, pp. 376-383, June 1988.
- [DaVe89] S. Damaskos and D. C. Verma, "Fast Establishment of Real-Time Channels", Tech. Rept. TR-89-056, International Computer Science Institute, Berkeley, October 1989.
- [Ferr89] D. Ferrari, "Real-Time Communication in Packet Switching Wide-Area Networks", Tech. Rept. TR-89-022, International Computer Science Institute, Berkeley, May 1989.

- [Ferr90] D. Ferrari, "Client Requirements for Real-Time Communication", Tech. Rept. TR-90-007, International Computer Science Institute, Berkeley, March 1990.
- [FeVe90] D. Ferrari and D. C. Verma, "A Scheme for Real-Time Channel Establishment in Wide-Area Networks", *IEEE J. Selected Areas Commun.*, vol. SAC-8, n. 3, pp. 368-379, April 1990.
- [Harr80] E. Harrington, "Voice/data integration using circuit-switched networks", *IEEE Trans. Comm.*, vol. COM-28, n. 6, pp. 781-793, June 1980.
- [Lein89] B. Leiner, "Critical Issues in High Bandwidth Networking", Internet RFC-1077, Nov. 1988.
- [OhON88] H. Ohnishi, T. Okada and K. Noguchi, "Flow control schemes and delay/loss trade-offs in ATM networks", *IEEE J. Selected Areas Commun.*, vol. SAC-6, n. 6, pp. 1609-1616, December 1988.
- [Schw89] H. Schwetman, "CSIM Reference Manual (Revision 12)", MCC Tech. Rept. No. ACA-ST-252-87, January 1989.
- [Zhan87] L. Zhang, "Designing a new architecture for packet switching communication networks", *IEEE Commun. Magaz.*, vol. 25, n. 9, pp. 5-12, September 1987.