# Real-time Communication in Multi-hop Networks

Dilip D. Kandlur
Kang G. Shin

Domenico Ferrari

Real-Time Computing Laboratory
Department of EE & CS
The University of Michigan
Ann Arbor, Michigan 48109

Computer Science Division
EE & CS Department
University of California
Berkeley, CA 94720

## Abstract

In this paper, we develop a scheme for providing predictable inter-process communication in real-time systems with (partially connected) point–to–point interconnection networks, which provides guarantees on the maximum delivery time for messages. This scheme is based on the concept of a real-time channel, a unidirectional connection between source and destination. A real-time channel has parameters which describe the performance requirements of the source–destination communication, e.g., from a sensor station to a control site. We concentrate on methods to compute guarantees for the delivery time of messages belonging to real-time channels. We also address problems associated with allocating buffers for these messages and develop a scheme which preserves delivery time guarantees.

## 1  Introduction

The problem of providing predictable inter-process communication is of great significance in real-time systems because unpredictable delays in the delivery of messages can affect the completion time of tasks participating in the message communication. If these tasks have deadlines associated with them, communication delays could cause the tasks to miss their deadlines. In this paper, we will address the issue of guaranteeing the delivery of messages with time constraints in multicomputers with point–to–point interconnection networks.

This problem of time-constrained communication has been studied by several researchers, since it plays an important role in video– and voice– data transmission over a data network. Recently, it has also been studied in the context of communication in embedded or real-time systems. These efforts have been directed mainly towards designing medium access protocols for multiple-access networks which consider time constraints on messages. The survey paper by Kurose *et al.* [1] discusses many of the proposed techniques. Most of these schemes can be classified as *best-effort* schemes, where the system tries to ensure that most

messages meet their deadlines, but it cannot give any guarantees about the delivery times. On the other hand, when the system has some information about the arrival pattern of messages, it can try to give guarantees about their delivery times. For example, Strosnider and Marchok [2] use a variation of the rate-monotonic scheduling algorithm to control access to a token-ring network. They assign priorities at design time to message sources based on the periodicity of message generation and they can check for the possibility of deadline overrun at that time. In their work, the time-constraint on periodic messages is implicitly assumed to be the beginning of the next period.

The network under consideration here is, however, not a multiple-access network, but a point–to–point interconnection structure which has potential for higher performance and reliability than bus/ring structures. In this case, the problem is more complicated than multiple–access networks, because we have to consider delivery time constraints across multiple stages in the network. In this type of network, there is only one source node for any network link, so the issue to be addressed is not one of access to the medium but that of message scheduling in the network nodes. Although it is possible that messages which have to traverse multiple links to reach the destination may suffer from the problem of higher latency, the latency can be made more *predictable* using message scheduling and network flow control. There has not been much work reported on the problem of providing time-constrained communication in a local-area point–to–point interconnection network. However, the work by researchers in the DASH and Tenet projects [3, 4] dealing with continuous–media communication in wide-area networks is closely related to the work reported here.

In this paper, we will be concentrating on issues related to the design of time-constrained communication in which the system provides *a priori* guarantees for message delivery. We will also briefly discuss how best-effort delivery can easily be accommodated into the design. Section 2 describes the computing environment and defines the problems to be addressed. In Section 3 we present a solution scheme which is general and works for messages of arbitrary size. Section 4 describes a buffer management scheme

which is consistent with the message handling scheme de-scibed in Section 3. We discuss other related work in Section 5 and draw our conclusions in Section 6.

## 2  Problem Statement

The intended application for this work is in embedded real-time systems, and thus, we make assumptions appropriate for that environment. The system consists of processor nodes connected by point–to–point dual simplex links. A node in the system can have several incoming and outgoing links connected to it, and these links can operate in parallel. The nodes run a common distributed real-time operating system which is responsible for control of the network. The operating system also synchronizes clocks and maintains a global time-base for the system, where the maximum skew between clocks on different nodes is very small compared to the end–to–end message delivery delay. A scheme for achieving this synchronization is presented in [5].

Communication between user-level entities in this system can either be connection-oriented or connectionless. Likewise, in the connection-oriented case, it can be either message stream or byte stream oriented. In the context of time-constrained communication, the communication model has to preserve message boundaries, and so a byte stream model is not suitable. Connection-oriented service is considered more suitable for applications which require guaranteed de-livery time for communication [3]. In order to make a guarantee about delivery time for a message, the system requires information about other message sources that can contend for resources with this message. It is therefore necessary to provide a mechanism for describing the characteristics of communication, so that resources can be reserved for real-time connections. The connection establishment procedure gives the service provider, in this case the distributed operating system, the opportunity to reserve resources for the connection, and for the user to specify its requirements. Therefore, the abstraction that we use for guaranteed time-constrained communication is one of connection-oriented sequenced messages, which we call a *real-time channel* or simply *channel*. In a bidirectional connection between a pair of user entities, the message generation characteristics may differ substantially for the two directions. Hence, it is preferable to restrict the real-time channel to unidirectional communication, and a bidirectional connection can be composed from a pair of channels.

The resources required by a channel include network bandwidth, buffer space, and message processing band-width. The operating system can make resource reservations based on the communication requirements of the user, as specified in the request for the connection. These requirements consist of the source and destination end-points, a description of the message generation process, and the desired end–to–end guarantee for the delivery time. Without a description of the message generation process, the service provider cannot compute the resource requirements and hence it cannot provide a guarantee.

It is perceived that a significant portion of the traffic which requires guarantees will be periodic communication, so the model for message generation is based on this perception. The message generation process is specified in terms of a linear bounded arrival process, a model which was first proposed by Cruz [6]. This model has also been adopted by Anderson and others [4] for continuous media applications, and some of the terminology given below is from [4]. The arrival process has the following parameters:

| | | |
|---|---|---|
| maximum message size | $S_{max}$ | (bytes) |
| maximum message rate | $R_{max}$ | (messages/second) |
| maximum burst size | $B_{max}$ | (messages) |

The model includes the restriction that, in any time interval of length $t$, the number of messages generated may not exceed $B_{max} + t \cdot R_{max}$, and that the length of each message may not exceed $S_{max}$. $R_{max}$ is a bound on the message generation rate and its reciprocal, $I_{min}$, is the minimum (logical) inter-arrival time between messages. The burst parameter $B_{max}$ puts a bound on the allowed short-term variation in the message generation rate, and partially determines the buffer space requirement of the channel. Message generation which is not periodic can be represented in this model using an estimate of the worst-case inter-arrival time and the average rate of generation.

The semantics of the end–to–end guarantee on delivery time are also based on this model. The *logical generation time*, $\ell(m)$, for a message $m$ can be defined as

$$\begin{aligned} \ell(m_0) &= t_0 \\ \ell(m_i) &= max\{\ell(m_{i-1}) + I_{min}, \ t_i\} \end{aligned}$$

where $t_i$ denotes the actual generation time of $m_i$. If $d$ is the end–to–end delay for the channel, the system *guarantees* that any message $m_i$ will be delivered to the destination node by time $\ell(m_i) + d$. In other words, when the inter-arrival time between messages is at least $I_{min}$, the system guarantees that each message in the channel incurs a delay of at most $d$ seconds. However, messages which arrive in a burst, where the inter-arrival time is less than $I_{min}$, may suffer a larger delay since the guarantee is given with respect to the logical arrival time. This possible increase in delay is a consequence of regulation: arrivals at each node have to be regulated in order to prevent burst arrivals on one channel from affecting the guaranteed delay of messages belonging to other channels.

This paper deals with problems connected with computing a guaranteed end–to–end delay for messages belonging to a channel, and the scheduling of messages to achieve this goal. These two problems are related because it is necessary to understand the scheduling environment in order to compute a guarantee. In the next section, we will first present a simple channel establishment procedure in which the problem of computing an end–to–end guaranteed delay is reduced to the problem of guaranteeing the worst-case delay for a single station/node.

1. Select a source–destination route for the channel.

2. Compute the worst-case delay for a message on each link on the route. In this computation, it is necessary to ensure that the new channel does not affect the guaranteed delivery times of existing channels. Also compute the buffer requirement for this channel.

3. Compute the sum of the worst-case delays which were determined in Step 2, and check whether it is less than the user-specified delay. This is the channel establishment test.

4. If the channel can be established, divide the user-specified delay among the links on the route based on their worst-case delay for the message. Adjust the buffer requirements based on these allocated delays.

Figure 1: **Channel Establishment Procedure**

## 3 Basic Solution Approach

There are two distinct phases in handling real-time channels: channel establishment and run-time message scheduling. The channel establishment phase is outlined in Figure 1, and begins with the selection of a route for the channel. There are two alternatives for routing packets in the network: dynamic routing and static routing. Dynamic routing offers the advantage that it can adapt to the network load and can reduce the average delivery delay for messages. However, it is very difficult to make any guarantees on message delivery for a channel based on dynamic routing, in which a message can use one of several alternate routes through the network. Therefore, we use static routing for messages belonging to real-time channels. We will assume that a route has been selected for the channel, possibly using a scheme like the one presented in [7], and concentrate on the computation of guarantees and on buffer management.

The worst-case delay for a message at a link depends upon the other channels which use the link and upon the scheduling algorithm. The link may also be used by messages without time constraints, but the effect of these messages on the worst-case delay is limited (i.e., restricted to a single packet delay if packet transmission cannot be aborted while in progress) because they would belong to a lower priority class. Hence, we can restrict our attention to messages belonging to real-time channels while calculating the worst-case delay. The scheduling environment that we encounter in message scheduling is one of independent, possibly periodic, message arrivals which have deadlines associated with them. The message deadline may be related to its period, but it is not necessarily the beginning of the next period. There are several approaches to scheduling these messages, which can be categorized as fixed priority or dynamic priority algorithms. For example, Earliest Due Date [8] is a dynamic priority algorithm, whereas rate monotonic scheduling [9] is a fixed priority algorithm. We will now discuss the problems associated with the use of these algorithms

for message scheduling and channel establishment.

**Deadline Scheduling:** The Earliest Due Date (EDD) algorithm schedules messages in the order of their deadlines, with higher priority given to messages which have closer deadlines. Liu and Layland [9] have shown that EDD is optimal for preemptive scheduling of periodic tasks when the task deadline is equal to its period, and that a feasible schedule exists whenever the total utilization is less than one. However, there is no similar result (based on utilization alone) available when the task deadlines are not related to their periods. The main drawback of EDD scheduling is that computation of guarantees is difficult, since the priority of a task depends upon the relative order of arrival of the tasks. A multiclass version of the EDD algorithm has been used for scheduling real-time messages by Ferrari and Verma [3], who also present sufficient conditions for the existence of feasible schedules. Their approach will be discussed in greater detail in Section 5.

**Fixed Priority Scheduling:** Scheduling decisions can be based on a fixed (static) priority scheduling algorithm where messages are processed and transmitted in the order of priority. For example, in rate-monotonic scheduling the priority assigned to a channel is related to the frequency of occurrence of messages on that channel. For any priority assignment scheme, if message arrivals on all the channels are assumed to be strictly periodic, we can determine whether the set of channels is schedulable. This is done by computing the worst-case response time for each message using a *critical time zone* analysis similar to the one used by Liu and Layland [9], and verifying that the worst-case response time is less than the delay assigned to that channel.

There are some problems with priority based scheduling and the computation of guarantees when we consider multiple stage service, like service for a message which has to traverse multiple links. The response time for the message varies depending upon the arrival time of other messages at a node. Therefore, even if messages are generated with a fixed inter-arrival time at the source, the inter-arrival time for the message at the next service point is not constant. Early arrivals are also possible due to burstiness in the message generation at the source. A message which arrives early at a node can cause a lower priority message to miss its deadline. Figure 2 shows the results of preemptive priority scheduling for two message streams, $M_1$ and $M_2$, with service requirements $C_1$ and $C_2$, where $M_1$ has higher priority than $M_2$. The figure shows how the early arrival of a message on $M_1$ can cause a message on $M_2$ to miss its deadline $d_2$. Note that a similar effect can be observed even when a non-preemptive discipline is used. One solution for this problem is to "hold back" the early arrival and not consider it for transmission until its scheduled arrival time. However, this scheme involves the setting and resetting of timers and is expensive to implement. Also, it implies that the message cannot make the best progress possible on lightly loaded links.
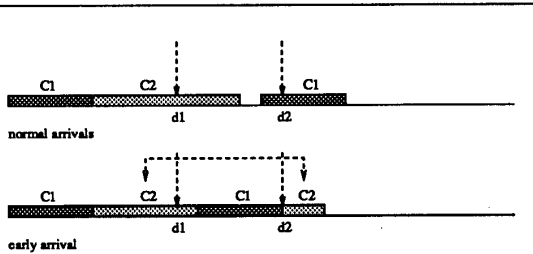
Figure 2: Effect of early arrivals.

## 3.1 Proposed Algorithm

The discussion on the problems of deadline and fixed priority scheduling in a multi-hop communication system suggests the use of a combination of deadline and fixed priority scheduling. Our channel establishment scheme uses fixed priority scheduling for computing the delay, but we use a form of EDD as the run-time message scheduling algorithm. When a channel is to be established, for each link on its route, we estimate the worst-case response time for a message based on fixed-priority scheduling. The details of this scheme have been omitted from this version due to space limitations. The response time analysis is based on an approximation of preemptive scheduling. Pure preemptive scheduling cannot be used in the context of message scheduling, because, if the transmission of a message is interrupted, the message is lost and has to be retransmitted. To achieve the benefits of preemptive scheduling, the message has to be split into packets so that message transmission can be interrupted at the end of a packet transmission, without loss. (This is analogous to allowing an interrupt at the end of an instruction execution.) Therefore we consider a message to be a set of one or more packets, where the packet size is bounded. Packet transmission is non-preemptive, but *message* transmission can be considered to be preemptive.

The priority assigned to the new channel at a link depends upon the characteristics of the other channels going through the link. The total response time, which is the sum of the response times over all the links on the route of the channel, is checked against the maximum permissible message delay and the channel can be established only if the latter is greater. In this case, the permissible message delay is split proportionally among the different links. While using this procedure, it is necessary to ensure that the new channel does not affect the guaranteed delivery times of existing channels. This is taken care of in the priority assignment algorithm.

## 3.2 Priority Assignment

The worst-case response time computation for a new channel at a link on the route requires an assignment of priorities for existing channels that use this link. This priority assignment problem can be defined formally as follows. Let $\{M_i = (C_i, p_i, d_i), i = 1, \ldots, k\}$ be the set of $k$ existing channels through a link $\ell$, where $C_i$ is the maximum service time required for messages of channel $M_i$ on this link, $p_i = I_{min}^i$ is the message inter-arrival time, and $d_i$ is the permissible delay which has been assigned to the channel for this link. The service time requirement $C_i$ is proportional to the product, $S_{max}^i . R_{max}^i$, of the maximum message size and the maximum message rate for the channel $M_i$. Given a new channel $M_{k+1} = (C_{k+1}, p_{k+1})$ to be established (the delay bound need not be specified on a per-link basis), we have to find a priority assignment for this augmented set of channels such that the response time $r_i'$ computed for a channel under this priority assignment satisfies the constraint $r_i' \leq d_i, i = 1, \ldots, k$. An assignment which satisfies this constraint is called a *feasible priority assignment*. Moreover, we would like to find a feasible priority assignment which will <u>minimize</u> the response time $r_{k+1}'$ for the new channel $M_{k+1}$. By finding the minimum feasible response time on each link, we can improve the chances of satisfying the channel establishment test.

Consider the procedure D_Order, shown in Figure 3, which assigns priorities and computes the response time for all channels, including the one to be established. This procedure works under the condition that, for all channels, $d_i \leq p_i$. That is, the worst-case delay at each link for any channel does not exceed its inter-arrival time. Note that the total end-to-end delay can exceed the inter-arrival time of the channel.

We can prove that, when $d_i \leq p_i, \forall i$, assignment procedure D_Order is optimal in the sense that the computed response time $r_{k+1}'$ is the <u>minimum possible</u> for any feasible priority assignment. Lemma 1 states that the priority assignment based on delays is optimal in the sense that if there is any feasible priority assignment for the channel set, there is a feasible priority assignment based on delays. This is in fact a generalization of the optimality result for rate-monotonic scheduling [9]. In their model, Liu and Layland assumed that the deadline for an instance of a periodic task is the release time of the next instance, that is, $d_i = p_i$. The rate-monotonic algorithm assigns priorities to periodic tasks based on their periodicity, with high priority assigned to tasks with high periodicity. The optimality result for this scheduling algorithm can be derived from Lemma 1 by substituting $d_i = p_i, \forall i$. It has been brought to our attention that a result similar to Lemma 1 has also been proved by Leung and Whitehead [10].

**Lemma 1:** Consider the set of channels $\{M_i = (C_i, p_i, d_i), i = 1, \ldots, k\}$ through a link. Assume that there exists a feasible fixed priority assignment for the channels such that the computed response time for each channel satisfies the constraint $r_i \leq d_i \leq p_i, \forall i$. Then, the priority assignment $P^D$, based on an increasing order of delays with high priority assigned to channels with small delays, is also a feasible priority assignment.

**Proof:** *omitted from this version.*

303

1. Arrange the channels in ascending order of their associated delay $d_i$.

2. Assign the highest priority to the new channel $M_{k+1}$. Assign priorities to the other channels based on this order, with high priority assigned to channels with small delays.

3. Compute the new (worst-case) response times $r_i'$ for the existing channels based on this priority assignment.

4. In the priority order, find the smallest position $q$ such that $r_i' \leq d_i$ for all channels with position greater than $q$.

5. Assign priority $q+1$ to the new channel and compute the response time $r_{k+1}'$.

---

**Figure 3: Assignment Procedure D_Order**

---

**Theorem 1:** Assignment procedure D_Order yields the smallest response time for the new channel, $M_{k+1}$.

**Proof:** Suppose there is some other feasible priority assignment $P$ which yields a smaller response time. If $r_{k+1}$ is the response time for $M_{k+1}$ under $P$, we can set $d_{k+1} = r_{k+1}$. Consider any channel $M_j$ which has lower priority than $M_{k+1}$. Then, $r_{k+1} < r_j \leq d_j$ and hence $d_{k+1} < d_j$. Now, by Lemma 1, the priority assignment $P^D$ based on the channel delays is also a feasible assignment. Under this assignment, $d_{k+1} < d_\ell$ implies that $M_{k+1}$ has priority over $M_\ell$. Therefore, any channel which has priority lower than $M_{k+1}$ under $P$ also has priority lower than $M_{k+1}$ under $P^D$. Also, a channel which has priority higher than $M_{k+1}$ under $P$ could have priority lower than $M_{k+1}$ under $P^D$. This means that the response time for $M_{k+1}$ under $P^D$ will be at most equal to the response time under $P$. Contradiction. ∎

In the last stage of the channel establishment procedure, delays will be assigned to the new channel, for each link on the source–destination route. For any link, the assigned delay for that link, $d_{as}$, has to be such that $d_{as} \geq r_{k+1}'$. It can be shown that if channels are allotted priorities based on this assigned delay, the resulting priority order is still feasible. The new priority order for the channels is not necessarily the one used to compute the response time $r_{k+1}'$ because there may exist a channel with $r_{k+1}' \leq d_i < d_{as}$.

## 3.3 Run-time Scheduling

A node in the system can have several incoming and outgoing links connected to it. These links can operate in parallel, so each outgoing link is considered as a separate entity for scheduling. Messages are composed of packets, and packets carry information about the message, and the channel, to which they belong. When a packet arrives at a node, or is generated in the node, it is dispatched to the appropriate outgoing link. Suppose the arriving packet belongs to the $i$th message on the channel $M_c$. All packets

belonging to the $i$th message would be assigned the same *logical arrival time*, $\ell_c(m_i)$, which is the logical arrival time for the message. For the source node $s$ of channel $M_c$, $\ell_{c,s}(m_i)$ is defined exactly in the same way as the logical generation time in Section 2, where $t_i$ denotes the generation time of message $m_i$.

$$\ell_{c,s}(m_0) = t_0$$
$$\ell_{c,s}(m_i) = max\{(\ell_{c,s}(m_{i-1}) + I_{min}^c), \ t_i\}.$$

For nodes other than the source node, the logical arrival time at the node is based on the logical arrival time of the message at the upstream node. Consider two adjacent nodes $a$ and $b$, sharing a link $(a,b)$ which is a part of the route for channel $M_c$. The logical arrival time for $m_i$ at node $b$, $\ell_{c,b}(m_i)$, is defined as

$$\ell_{c,b}(m_i) = \ell_{c,a}(m_i) + d_{c,a}$$

where $d_{c,a}$ is the worst-case delay for messages on channel $M_c$ at node $a$. It can be seen that the logical arrival time of a message at any node is ultimately based on its logical arrival time at the source node. This definition of the logical arrival time is feasible because the nodes in the system have synchronized clocks and the maximum skew between clocks on different nodes is small compared to the message delivery delays [5].

The logical arrival time that is assigned to messages is used by the message scheduler, which uses a variation of the multiclass EDD algorithm. In the following description, the second subscript has been dropped from the notation for the logical arrival time, since we are dealing with a single node. The scheduler maintains three queues for each outgoing link, corresponding to three service classes.

**Queue 1** Packets belonging to real-time channels with $\ell_c(m_i) \leq$ current_time, arranged in the order of increasing deadlines.

**Queue 2** Other packets arranged in the order of increasing deadlines.

**Queue 3** Packets belonging to real-time channels with $\ell_c(m_i) >$ current_time, arranged in the order of increasing logical arrival time.

Queue 1 and Queue 3 contain packets which belong to real-time channels, while Queue 2 contains all other types of packets. Queue 1 contains *current* real-time packets, which have to be scheduled in the order of their deadline, hence it is organized by increasing packet deadlines. Current real-time packets are those packets whose logical arrival time is less than the current clock time at the node. Queue 2 contains packets for which no guarantees are given, but which can have deadlines associated with them. For example, packets belonging to messages which request "best-effort" type service would fall in this category. The service provided for this class of packets is improved by giving them priority over real-time packets which are not *current*. Since they are not pertinent to the current discussion, we will not

1. find $\ell_c(m_i)$, the logical arrival time of the message to which this packet belongs.

2. set $\ell t(\text{packet}) = \ell_c(m_i)$.

3. set the deadline for this packet to $\ell_c(m_i) + d_c$.

4. if $(\ell t(\text{packet}) \le \text{current\_time})$
       insert packet into Queue 1
   else
       insert packet into Queue 3

5. invoke the dispatcher.

---

**Figure 4: Processing on Packet Arrival**

---

elaborate on the treatment of this class any further. Packets in Queue 3 are those which have arrived early, either because of burstiness in the message generation or because they encountered delays which were smaller than the budgeted worst-case delays at some upstream service stations. These packets are stored in the order of their logical arrival time because they have to be transferred to Queue 1 as they become current.

The actions taken when a real-time packet arrives are shown in Figure 4, whereas non real-time packets are simply inserted into Queue 2. The logical arrival time is determined based on the channel to which the packet belongs and the sequence number of the message of which it is a part. This is the logical time, $\ell t()$, of a packet. The deadline for the packet is set to $\ell_c(m_i) + d_c$, where $d_c$ is the worst-case delay guaranteed for channel $M_c$ at this node. Since all packets in a message have the same $\ell_c(m_i)$ and the same $d_c$, they will all have the same deadline. The packet is then inserted into Queue 1 or Queue 3, depending upon whether it is current or early. Lastly, the dispatcher is invoked.

The dispatcher, which is shown in Figure 5, first checks whether any real-time packets in Queue 3 have become current and transfers such packets to Queue 1. If the link is idle, it examines the queues in the order of priority looking for a packet to transmit. Packets in Queue 3 are considered for transmission only if their logical time is within the *horizon*, $H_k$, for the link. The horizon is link-dependent and is used for flow control, as explained in Section 4. These packets are scheduled in the order of logical time, primarily because they are queued in that order. Also, this ordering makes it easy to identify packets which are within the horizon (and can be considered for transmission). The dispatcher is also invoked upon completion of transmission of a packet on the link. In some situations, it is possible that Queues 1 and 2 are empty and Queue 3 only contains packets which are ineligible for transmission. In such cases, a timer can be used to trigger the dispatcher at the appropriate future instant.

## Remarks on Correctness

We now have to establish that this run-time scheduling scheme will conform with the guarantees computed using

0.  Examine Queue 3. Transfer those packets which have $\ell t(\text{packet}) < \text{current\_time}$, to Queue 1.
1.  if (link idle)
2.      if (Queue 1 nonempty)
3.          start transmission of head(Queue 1)
4.      else if (Queue 2 nonempty)
5.          start transmission of head(Queue 2)
6.      else if (Queue 3 nonempty)
7.          if $(\ell t(\text{head}(\text{Queue 3})) < \text{current\_time} + H_k)$
8.              start transmission of head(Queue 3)
9.          end
10.     end
11. end

---

**Figure 5: Dispatcher**

---

fixed message priorities. We approach this in two stages. In the first stage, we show that if we use this scheduling scheme on a channel set (which is schedulable) that satisfies the inter-arrival time constraints, then no message will miss its deadline. In the second stage we consider the effects of early arrivals and show how our scheme can accommodate them.

**Stage 1:** The system times for channels have been computed based on priority scheduling, taking into account the maximum rate of arrival of messages; and these times satisfy the local deadlines for each channel. Hence, when arrivals conform to the inter-arrival time constraints, a feasible schedule based on priorities exists. Dertouzos [8] has shown that this implies the existence of a feasible EDD schedule for this message set, which is based on the deadlines assigned to the channels.

**Stage 2:** Consider a message $m_j$ belonging to channel $M_c$, which arrives at the source node before its logical arrival time, that is, $t_{arr,m_j} < \ell_c(m_j) = \ell_c(m_{j-1}) + I^c_{min}$. This message would be assigned a deadline, $D_j = \ell_c(m_j) + d_c$, corresponding to its logical arrival time. By the argument given above, we can assert that, if the message did come in at its logical arrival time, a feasible EDD schedule exists by which all messages will meet their deadline. Since scheduling is based on deadlines, $m_j$ can affect only those messages which have deadlines beyond $D_j$. However, these messages would be affected to the same extent even if $m_j$ arrived at its logical arrival time, since $D_j$ is assigned based on the logical arrival time. From this, it follows that a feasible EDD schedule exists even in this case. The same argument can be applied even when multiple messages arrive prior to their logical arrival time.

If a message $m_i$ arrives late at the source node, that is, $t_{arr,m_i} > \ell_c(m_{i-1}) + I^c_{min}$, it will be assigned a logical arrival time $\ell_c(m_i) = t_{arr,m_i}$ and a deadline, $D_i = \ell_c(m_i) + d_c$. Also, future message arrivals on this channel will be assigned logical times based on $\ell_c(m_i)$. This mes-

sage cannot increase the maximum system time for any other messages, since the system times for these messages were computed by considering the maximum possible arrival rate for messages on channel $M_c$. Hence, a feasible priority-based schedule exists in this case, and consequently, a feasible EDD schedule also exists. The logical time assigned to messages at intermediate nodes is based on the logical time at the source node, and similar arguments apply.

## 4 Buffer Management

Buffer space has to be reserved for real-time channels at the source, destination, and intermediate nodes in order to prevent buffer overruns and consequent loss of messages. During channel establishment, the user has to specify the burstiness in message generation, $B_{max}$, as part of the channel specification. This burstiness determines the buffer requirement at the source node of the channel. If $d_{c,s}$ is the worst-case delay guaranteed for channel $M_c$ at the source node, the buffer requirement can be expressed as $(B_{max}^c + d_{c,s} \cdot R_{max}^c)S_{max}^c$, where the $d_{c,s} \cdot R_{max}^c$ term accounts for new arrivals during the system (response) time of a message on the channel. It is necessary for the source node to provide this buffer space because the client can legitimately produce messages at this rate, and, if buffer space is not provided, messages could be lost due to non-availability of buffers. Also, the source node is responsible for making sure that message generation for $M_c$ adheres to the channel specification. It can refuse to accept messages belonging to $M_c$ when the actual generation rate exceeds the specified bounds.

Intermediate nodes also have to provide buffer space for $M_c$, but this need not depend upon $B_{max}^c$ if a flow-control mechanism is employed between nodes. The buffer reservation scheme and the flow-control mechanism are intimately related, since flow-control can be used to regulate the burstiness of message arrivals at intermediate nodes. The flow-control mechanism considered here operates by holding back some of the messages which arrive before their logical arrival time. If flow-control were not used, messages could zip through lightly loaded nodes and back up at heavily loaded nodes, thereby causing buffer space problems at the heavily loaded nodes. An extreme case of flow-control is when all messages which arrive early are held back, that is, a message is considered for transmission only when $\ell_c(message) \leq$ current_time. In this extreme case, an intermediate node, $k$, need only provide buffer space proportional to $d_{c,k} \cdot R_{max}^c \cdot S_{max}^c$ and rely on flow-control to eliminate the burstiness. However, this strict regulation can result in unnecessary delays for messages when the network is lightly loaded. By allocating more buffer space, we make it possible for messages to go quickly through lightly loaded nodes.

The logical arrival time assigned to a message depends upon the burstiness in the arrival process. Hence, the bound on burstiness can also be interpreted as a *horizon*

for the logical arrival time. Burstiness can be measured by the equation,

$$B = (\ell_c(\text{message}) - \text{current\_time})/I_{min}^c$$

which can be rewritten as

$$\ell_c(\text{message}) = \text{current\_time} + B \cdot I_{min}^c.$$

A bound on the burstiness, $B \leq B_m$, can therefore be interpreted as $\ell_c(\text{message}) \leq \text{current\_time} + B_m \cdot I_{min}^c$. Packets belonging to messages with logical arrival time greater than this horizon may be discarded by the node. In other words, the horizon, $B_m \cdot I_{min}^c$, determines the buffer requirement for a channel at an intermediate node.

A node can compute the horizon for a link, based on the available buffer memory, whenever a new channel is to be established with that link. This horizon is passed back to the upstream node, since that node is responsible for flow-control on the link. The upstream node has to make sure that it transmits forward only those messages which have logical time, at that node, within the horizon (packets which belong to a single message all have the same logical time). It is possible to use a different horizon for each channel, but this can create problems for the run-time scheduler because it is then forced to examine all the packets in Queue 3 individually to check their eligibility for transmission. Use of a single horizon means less flexibility in buffer management, but it is more suitable for run-time scheduling. The buffer space required for a channel is computed as follows.

Consider two adjacent nodes $a$ and $b$, sharing a link $(a, b)$ which has a horizon $H_{a,b}$. Consider a message $p$ arriving at node $b$ on channel $M_c$, using link $(a, b)$. Since node $a$ uses $H_{a,b}$ for flow-control on link $(a, b)$, we can deduce that the logical time for the message at node $a$ must be $\ell_{c,a}(p) \leq \text{current\_time} + H_{a,b}$. The logical time assigned to the message at node $b$ is $\ell_{c,b}(p) = \ell_{c,a}(p) + d_{c,a}$, where $d_{c,a}$ is the assigned worst-case delay for $M_c$ at node $a$. Node $b$ will not contain any messages belonging to $M_c$ which have $\ell_{c,b}(m) < \text{current\_time} - d_{c,b}$, where $d_{c,b}$ is the worst-case delay assigned to the channel at this node. This is because messages on $M_c$ are guaranteed to leave $b$ within $d_{c,b}$ time units from their logical arrival time. Hence, node $b$ can only contain messages belonging to $M_c$ which have logical time in the range

$$\text{current\_time} - d_{c,b} \leq \ell_{c,b}(p) \leq \text{current\_time} + H_{a,b} + d_{c,a}.$$

The buffer space required to hold these packets is then given by:

$$\lceil (H_{a,b} + d_{c,a} + d_{c,b})/I_{min}^c \rceil \cdot S_{max}^c.$$

From this equation, it is clear that the minimum buffer space required for a channel is $S_{max} \cdot (d_{prev} + d_{node})/I_{min}$. The overall buffer space requirement of a node is the sum of the buffer space required for all channels going through the node. When a new channel is being established, the

306

node first tries to accommodate the new channel by computing the buffer space based on the existing horizon. If the available buffer space is less than the requirement, the node then reduces the horizon for that link. Reducing the horizon for a link has the effect of reducing the buffer space requirement of *all* channels using the link, so the amount of buffer space available increases.

## 5 Related Work

The concept of a unidirectional real-time channel was developed by the researchers of the DASH project. The real-time channel considered in this paper is similar to the deterministic real-time channel of [3]. That paper also considers two other types of channels: statistical, where messages are given a guarantee of successful delivery with a certain probability, and best-effort (with no guarantees). We do not consider statistical channels mainly because real-time applications require hard guarantees. Comer and Yavatkar have developed an abstraction called *flows* [11], which is similar to a unidirectional channel, but they do not give guarantees for the delivery time of messages. Flows can be identified with best-effort channels.

The parameters used for the description of the message arrival process are based on the linear bounded arrival process model of Cruz [6], with small modifications. Although the outline of our channel establishment scheme is similar to the one presented in [3], it is different in several respects. We have examined the problem of guarantee computation in depth, and developed a scheme for computing worst-case response times based on priorities. This corresponds to the delay bound test of [3]. It can be shown that the delay bound test of [3] is a special case of the priority-based response time computation scheme presented here. Moreover, we have presented an integrated solution to the problem of buffer reservation and flow control, based on the *horizon* model.

The problem of time-constrained communication in a multi-hop network has also been addressed by Cidon *et al.* in the PARIS system [12]. Their approach is based on providing limited buffer space in the switching nodes and using FIFO scheduling to obtain an upper bound on the network delay in each node. Limited buffering can result in packet loss, and so their efforts are aimed at maximizing throughput, while keeping the probability of packet loss below a certain level.

## 6 Summary

In this paper, we have identified and solved problems related to the establishment and management of real-time channels. We have presented algorithms for computing the worst-case delay for messages, and for scheduling these messages. The computation was based on priority scheduling, and it was shown to be optimal under certain conditions. We also presented mechanisms for buffer allocation and flow control suitable for real-time channels, which would preserve the guarantees on delivery time.

# References

[1] J. F. Kurose, M. Schwartz, and Y. Yemini, "Multiple-access protocols and time-constrained communication," *ACM Computing Surveys*, vol. 16, pp. 43–70, Mar. 1984.

[2] J. K. Strosnider and T. E. Marchok, "Responsive, deterministic IEEE 802.5 token ring scheduling," *Journal of Real-Time Systems*, vol. 1, pp. 133–158, Sept. 1989.

[3] D. Ferrari and D. C. Verma, "A scheme for real-time channel establishment in wide-area networks," *IEEE Journal on Selected Areas in Communications*, vol. SAC-8, pp. 368–379, Apr. 1990.

[4] D. P. Anderson, S. Y. Tzou, R. Wahbe, R. Govindan, and M. Andrews, "Support for continuous media in the DASH system," in *Proc. 10th Intl. Conf. on Distributed Computing Systems*, pp. 54–61, IEEE, May 1990.

[5] P. Ramanathan, D. D. Kandlur, and K. G. Shin, "Hardware assisted software clock synchronization for homogeneous distributed systems," *IEEE Trans. Computers*, vol. C-39, pp. 514–524, Apr. 1990.

[6] R. L. Cruz, *A Calculus for Network Delay and a Note on Topologies of Interconnection Networks*. PhD thesis, University of Illinois at Urbana-Champaign, July 1987. available as technical report UILU–ENG–87–2246.

[7] D. D. Kandlur and K. G. Shin, "Traffic routing for networks with virtual cut-through capability," in *Proc. 10th Intl. Conf. on Distributed Computing Systems*, pp. 398–405, IEEE, May 1990.

[8] M. L. Dertouzos, "Control robotics: The procedural control of physical processes," in *Proc. IFIP Congress*, pp. 807–813, 1974.

[9] C. L. Liu and J. W. Layland, "Scheduling algorithms for multiprogramming in a hard-real-time environment," *Journal of the ACM*, vol. 20, pp. 46–61, Jan. 1973.

[10] J. Y.-T. Leung and J. Whitehead, "On the complexity of fixed-priority scheduling of periodic, real-time tasks," *Performance Evaluation*, vol. 2, pp. 237–250, Dec. 1982.

[11] D. E. Comer and R. Yavatkar, "FLOWS: Performance guarantees in best effort delivery systems," Technical Report CSD-TR-791, Computer Science Department, Purdue University, July 1988.

[12] I. Cidon and I. S. Gopal, "PARIS: An approach to integrated high-speed private networks," *Intl. Journal of Digital and Analog Cabled Systems*, vol. 1, pp. 77–86, Apr. 1988.