

Adaptive Real-Time Group Multicast

Mario Baldi *

Yoram Ofek †

Bülent Yener ‡

Abstract

This paper shows how to provide an *adaptive real-time group multicast* (many-to-many) communication service. Such a service can be used by applications, like audio/video tele-conferencing, that require low loss, and bounded delay and jitter. In order to meet deterministic quality of service (QoS) requirements of a real-time group multicast, some communication resources are reserved. In this work we show (i) how bandwidth is reserved for each group, and (ii) how an active user in a multicast group can *dynamically* share, in an efficient and fair manner, the bandwidth allocated to its group.

Quality of service support for a real-time multicast group is based on *time driven priority* [3]. In this scheme the time is divided into *time frames* of fixed duration and all the time frames are aligned by using a global time reference which can be obtained from GPS (global positioning system). Bandwidth is allocated to a multicast group as a whole, rather than individually to each user. The allocation is done by reserving time intervals within time frames in some periodic fashion.

This sort of allocation raises two problems that are studied in this paper: (1) *Scheduling*: how time intervals are reserved to each multicast group, and (2) *Adaptive sharing*: how the participants dynamically share the time intervals that have been reserved for their multicast group. The proposed approach is based on embedding multiple *virtual rings*, one for each multicast group. By using the virtual rings it is simple to route messages to all the participants, while minimizing the bound on the buffer sizes and queuing delays.

Key words: multicast, real-time.

1 Introduction

Group or many-to-many real-time multicast is a demanding service. It is usually exploited by applications which require guaranteed (sometimes large) bandwidth with small loss, and bounded end-to-end delay and jitter^{**}. Given the finite network resources, applications should be performed in an *efficient* and *adaptive* manner. Efficient means that the communication links can operate near full utilization.

The number and the identity of transmitting nodes in a

real-time multicast group changes over time. This has two implications: (i) the resources should be allocated to the group and not individually to each participants, and (ii) the active nodes (i.e., the source nodes) should *adapt* their transmission rate for dynamic sharing of the bandwidth allocated to their multicast group.

The solution presented in this paper exploits a combination of *space* (virtual ring embedding) and *time* (packet forwarding driven by time). Using virtual rings for many-to-many communications is simple and “natural”, since a message sent from one member of the group travels around the ring until it comes back to the sender. We simply exploit this property by embedding multiple virtual rings, one for each multicast group. In a previous work it has been shown that the virtual ring can be instrumental in designing a *reliable multicast protocol* from *bursty sources* [4]. The ring circular property is used to adapt, in real-time, the share of capacity each *active node* obtains. It is done in a fair manner according to the number of nodes that are trying to transmit to the group at any given time.

An external timing information which comes, for example, from the global positioning system (GPS) [1], is used to control the forwarding of packets inside the network. This real-time forwarding is used for:

- (i) Minimizing the *queuing delay* bound in the virtual ring, as we discuss in Section 2.
- (ii) Sharing the capacity of a link among multiple virtual rings traversing the link in the same direction, as described in Section 3.
- (iii) Adaptive sharing of bandwidth to ensure that each active node in a virtual ring gets a fair share of the capacity, as discussed in Section 4.

2 Principles of Operation

2.1 Network Model

In this work we consider a general topology network with bidirectional links. A virtual ring (VR) is associated with each group multicast by constructing cycles that includes each participant at least once [4]. For example, Figure 1 shows two VRs (continuous and dashed lines) embedded onto a arbitrary topology network (bold lines). In this example, the two VRs share link (i,j) that will carry packets for both the VRs.

A virtual ring has three types of nodes: (i) active participants (ii) passive participants (e.g., listeners in a video conferencing), and (iii) the transit nodes that are not the members of the multicast group but they needed to construct the VR (shaded circles in Figure 1).

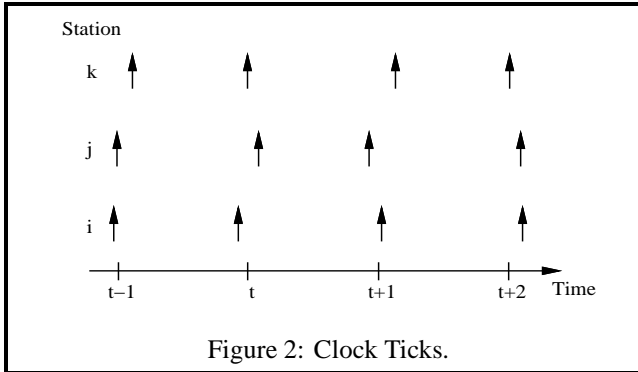
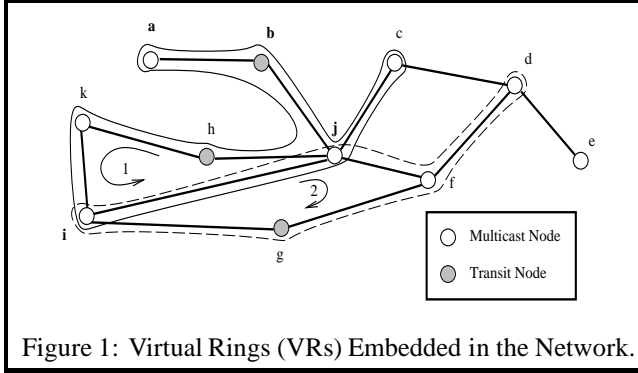
Although, some nodes are traversed more than once by the packets pertain to the same multicast group (e.g., f and j in Figure 1), these packets will be received only the first time

*Dip. di Automatica e Informatica, Politecnico di Torino, Torino, Italy 10129, baldi@leonardo.polito.it. This work was done in part while the author was visiting IBM T. J. Watson Research Center.

†IBM T.J. Watson Research Center, P.O.Box 218, Yorktown Heights, NY 10598, (914) 945-3171, ofek@watson.ibm.com.

‡Networking and Communication Lab. Department of Computer and Information Sciences, New Jersey Institute of Technology, Newark, NJ 07102-1982, (201) 596-2666, yener@cis.njit.edu. This work was supported in part by the Office of Naval Research under Grant N00014-96-1-0825.

**The jitter is defined as the difference between maximum and minimum delay.



they reach to these nodes.

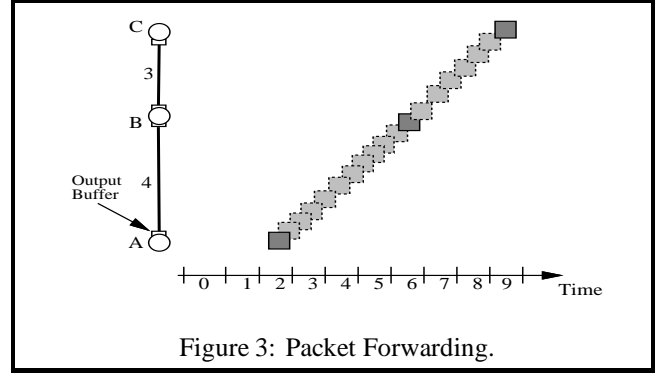
2.2 Network Timing

All the nodes are provided with a global timing structure. If, for example, GPS is used for this purpose, each node in the network receives clock ticks from a GPS receiver and uses them to keep its internal clock synchronized. As shown in Figure 2, there is a period of time in which the clocks of all the nodes have the same time value.

The time between two successive ticks is called a *Time Frame* (TF) and has a typical duration of $125 \mu s$ or 8KHz. Thanks to this global timing structure all nodes have the same knowledge of the *temporal value* of the current TF.

2.3 Packet Forwarding

Nodes forward packets according to the *pseudo-isochronous* or *RISC-like* forwarding principle proposed in [2, 3]. The basic idea is that packet hopping in the network is a function of TFs. Precisely, a packet received by a node in TF t , is forwarded by the node so that it reaches the next node along the path in TF $t+i$ (for $i \geq 1$). Note that in wide area internetworking environments i may be large due to the processing and propagation delays. Furthermore, this scheme assumes that the delay between the output port of a node and the output port of the next node along the path is a constant integer number of TFs. If this does not fit actual delays, receiving nodes can introduce extra delay by buffering packets and moving them to the output port at the predefined



instant.

Consider as an example the simple topology depicted in Figure 3; links are labeled with the delay (in TF units) between the output port of the node transmitting on the link, and the output port of the receiving node. Node A has a packet to send on its output port (buffer) at TF 2. This packet will arrive at node B and will be ready to be sent from B's output port at TF 6; the destination node C receives the packet at TF 9. The propagation delay from the source A to the destination C is fixed; the jitter is bound by two TFs and is due to the uncertainty on when node A's transmission begins, during TF 2, and when node C's reception ends, during TF 9.

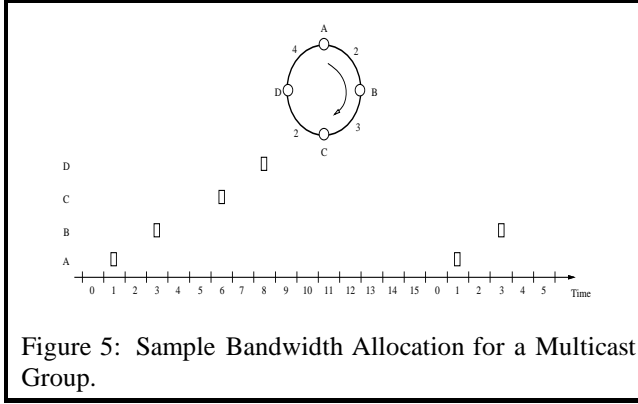
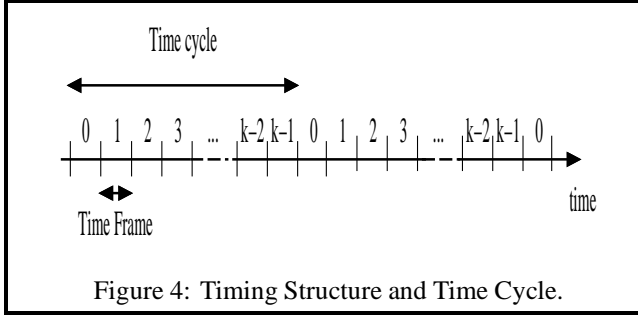
The packet forwarding on a VR follows the same principle such that a packet is forwarded along the VR and after a fixed number of TFs (say n_r) it is received back by its source node. Note that the integer n_r gives the *length* of the VR in terms of time frames.

2.4 Bandwidth Allocation

Some bandwidth is allocated to each multicast group and is shared by the members of the group that are ready to transmit. Section 3 shows how bandwidth allocation is performed while Section 4 shows how bandwidth is dynamically shared among the active nodes of the group.

We assume that a node in the network can forward from its output port up to $125 \mu s \cdot B M b/s$ bits during each TF. For example, if the data rate of a link is 100 Mb/s then 12.5 Kb can be transmitted in every TF. The source of the bits can be from this node, from its multicast group, or any other multicast group sharing this link. Note that $125 \mu s \cdot B M b/s$ is also the output buffer requirement on every port. This establish a *small bound* on the queueing delay of every node along the virtual ring. If the many-to-many group multicast would be done over a tree, the propagation delay would be shorter but the queueing delay bound would be much larger. This is significant, since the delay bound is a critical parameter for real time applications.

The bandwidth is allocated on a link for a VR (i.e., multicast group) by reserving time intervals in some TFs in a periodic manner. (Unused reserved bandwidth can be used for "best effort" traffic). For example, allocating a 64 Kb/s



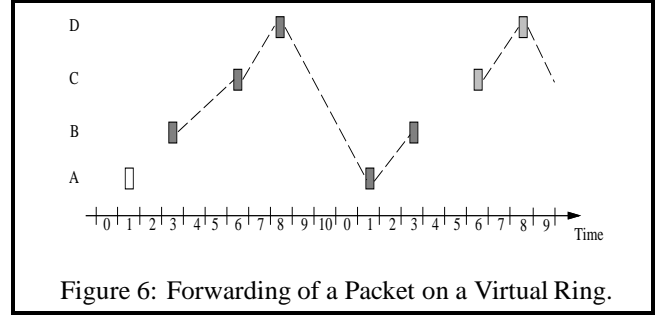
for a multicast group means the transmission of 1 byte^{††} (i.e., $0.064 \cdot 125 = 8$ bits) during each TF (in each virtual node). A packet size of one byte is obviously not practical. Thus, we have to use a packet size that keeps low the percentage overhead of the packet headers. For example for video-conferencing over IP networks, a 500 bytes payload is a reasonable size in order to keep the overhead due to IP, UDP, and RTP headers under 10%. Therefore, in our example, the 64 Kb/s can be allocated to a multicast group by reserving 500 bytes every 500 TFs.

TFs are grouped into *Time Cycles* (TCs) of fixed length k and are enumerated from 0 to $k - 1$, as shown in Figure 4. Allocation of bandwidth to a multicast group is done by reserving a fraction of the bits that a node in the VR can send during some of the k TFs in the TC.

Figure 5 shows a sample reservation of one packet per TC for a VR. In the figure a rectangle represents the bits reserved during the TF on the virtual ring nodes (i.e., A, B, C and D). Each node on the ring sends or forwards packets of its multicast group during the TFs in which transmission time for packets have been reserved, i.e., it uses always the same TFs in each TC for transmitting packets pertaining to its multicast group.

Assume for simplicity that packets are fixed size of s bytes. If we reserve a packet worth of bytes in p_r TFs for a multicast group r , then the bandwidth allocated to the multicast group

^{††}In the following examples when dealing with bandwidth we consider the gross bandwidth allocated to a multicast group, not the net bandwidth offered to applications. I.e., the overheads due to protocol encapsulation are not taken into account. This is not a limitation because these overheads can be taken into account by simply having the gross allocated bandwidth larger.



is given by

$$MB_r = \frac{p_r \cdot s \cdot 8}{k \cdot 125} \text{ Mb/s.} \quad (1)$$

Bandwidth is allocated in multiples of a base quantity which is given by the reservation of 1 packet per TC. In this work we assume $s = 500$ bytes and, as we want the base bandwidth allocation unit to be 64 Kb/s, we fix the number of TFs per TC (length of the TC) to

$$k = \frac{500 \cdot 8}{64 \cdot 10^3 \cdot 125 \cdot 10^{-6}} = 500.$$

Without loss of generality, we assume that only one packet is sent in each TF. As a consequence, bandwidth is allocated by giving to each virtual node on the ring exclusive usage of outgoing links during the whole reserved TFs. Such TFs are said to have been reserved on the node for a specific multicast group or VR.

2.5 Forwarding on a Virtual Ring

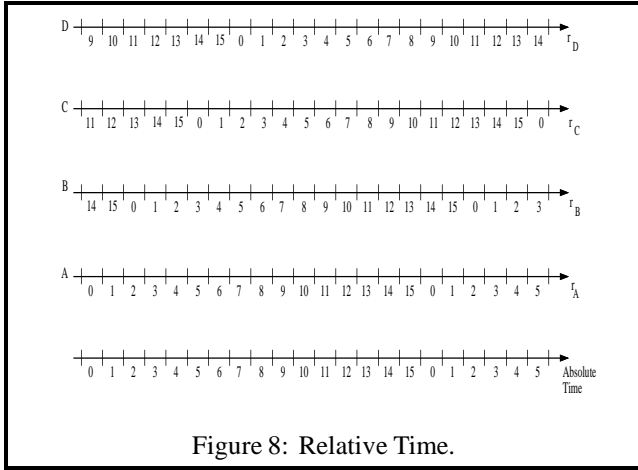
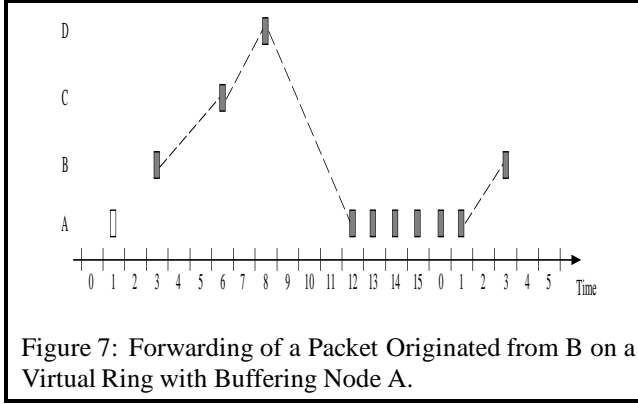
As an example, consider now the packet forwarding in the sample VR topology with the bandwidth allocation shown in Figure 5. If node B has a packet to transmit to the multicast group, it waits until TF 3 and then sends the packet which reaches the output port of C in TF 6, as shown in Figure 6. The packet travels on the VR and gets back to node B by TF 3 of the following TC. If node B at this point does not transmit any packet, C will be able to transmit a packet in TF 6 if it is not used for forwarding a packet from some other node.

Figure 7 shows the forwarding of packets on the same VR when $k = 16$. Virtual node A receives B's packet by TF 12, but it has no TF reserved until TF 1 of the following TC. Thus, node A buffers B's packet until the next TF reserved for its VR. Node A is named the *buffering node* of this VR. Each VR must have a buffering node in order to *adapt the length of the VR (virtual ring) to the length of the TC (time cycle)*. (Note that due to the presence of a buffering node a node does not necessarily receive its own packet at the same TF, in the next TC, in which it has transmitted it.)

2.6 Relative Time

When the timing structure is applied to a ring, the concept of *Relative Identifier* (RID) of TFs is useful for two purposes:

1. To handle the reservation of TFs in the nodes (Section 3) and



2. To dynamically share the allocated TFs among the participants to the multicast group (Section 4).

Each node j keeps its *relative time* by assigning a RID to TF i through the expression:

$$r_i(j) = (i - d_{b,j} + k) \bmod k \quad (2)$$

where $r_i(j)$ is the RID given to TF i by node j and $d_{b,j}$ is the delay, expressed in number of TFs, experienced by a packet traveling from node b chosen as the *base* for the relative time, to node j . The base node has the property that

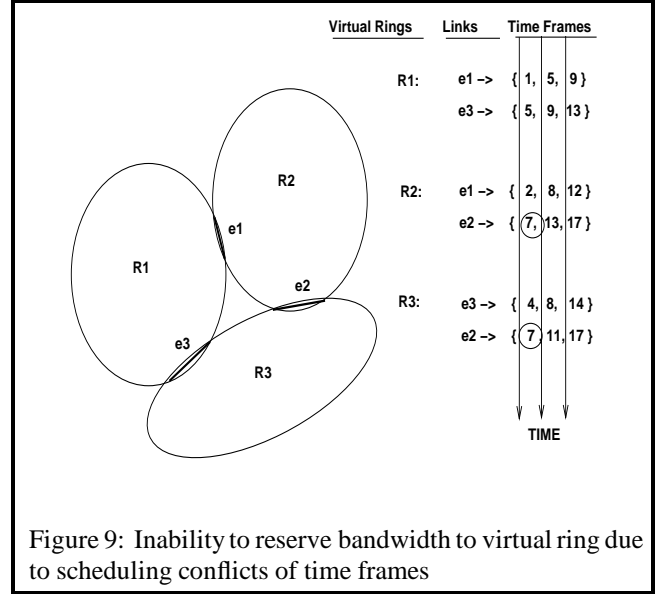
$$r_i(b) = i, \forall i \text{ } 0 \leq i \leq k - 1$$

, i.e., the relative time in the base node equals the absolute time.

Figure 8 shows the relative time for the nodes of the VR shown in Figure 5 when node A is the base node. By comparing the two figures, it is easy to note that the TF reserved to the VR has RID 1 in each node. In general, the relative time has the useful property that each node of a VR identifies the TFs reserved to the multicast group through the same RID.

3 Scheduling

In this Section we consider how to allocate resources to a new multicast group M_r with bandwidth request MB_r . As explained in section 2.4, the bandwidth request is mapped to a number (p_r) of TFs needed for the new multicast group.



Let VR_r be the VR associated with this multicast group and let E be the set of links included in this VR. Since the links of a synchronized network are shared by various types of traffic, some of the TFs (say f) on a link e in E may already be reserved.

It is necessary (*not sufficient*) to reserve p_r unused (available) TFs on each link in E , in order to guarantee MB_r bandwidth for this multicast group. Having p_r available TFs on each link is not sufficient because a TF (on the same link) may need to be allocated to two or more different VRs.

In other words, it is possible that there is available bandwidth (i.e., time frames) for a new multicast group (or to increase the allocation of an existing group), but there is no available *schedule* to assign time frames (of fraction of time frames) such that a time interval in a time frame is *not* allocated to more than one virtual ring, see an example in Figure 9. This problem is like the Pigeon Hole Problem where there are enough available TFs on each link, but there is no schedule that can fit them in the right linear (sequential) order, on each of the VRs. The impact of such scheduling conflict is that some of the link bandwidth cannot be allocated to the multicast groups (however, note that this bandwidth can still be used by non-reserved traffic).

The scheduling problem is examined in a hierarchical way. First, we consider a single link used by the multicast group and define the necessary condition for *feasibility* of a link. Next, we consider the whole VR and define *feasible ring scheduling*.

A link of a VR is *feasible* if it has at least p_r available TFs such that p_r satisfies the equation (1). Note that if there is even one link in E with less available TFs than p_r (i.e., the link is highly loaded), then this multicast group cannot be allocated the required bandwidth. Next, we introduce the following definitions to be used in this section:

Definition 1 –Feasible Schedule: Given a virtual ring

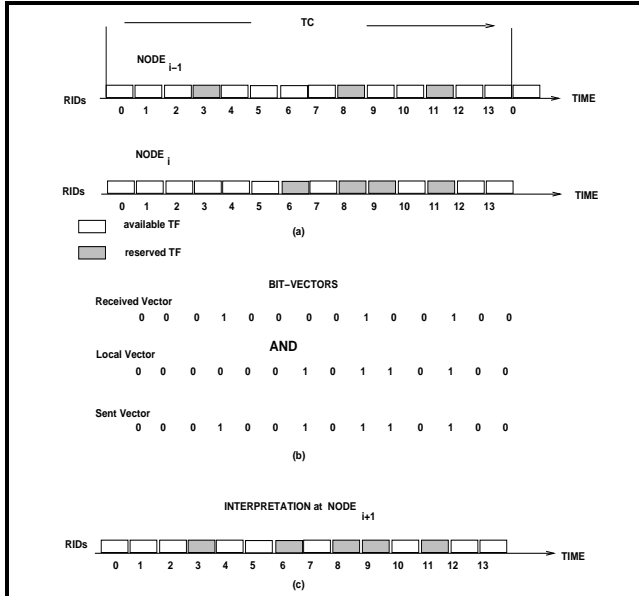


Figure 10: Example: Processing and Interpretation of Bit-Vectors to Determine a Feasible Schedule.

VR_r , a schedule σ is called *feasible schedule* for VR_r if σ chooses exactly p_r (available) TFs on each link such that each TF has the same RID at each node.

Definition 2 –Feasible Ring: A virtual ring VR_r is called *feasible* if there exists a feasible schedule σ for VR_r .

Definition 3 –Feasible Network: Given a network and a set of multicast groups on it, the network is called *feasible*, if there exists a feasible scheduling for each one of the VRs.

Proposition 1 In a feasible network, no TF on a link is used by two different rings.

3.1 Feasible Ring-Scheduling Algorithm

In this section we present a distributed algorithm to find a feasible assignment of TFs to a VR. The algorithm is initiated by a predetermined origin node (e.g., the base node) and executed by each node in a VR. In Figure 10 the steps of the algorithm are shown.

In the algorithm a node first assigns a RID to each TF by using equation (2). Next, a bit-vector BV_i of size k is constructed at node i . In BV_i a bit is set to 1 (see Figure 10.a), if the corresponding relative TF is not available (i.e., it has been previously reserved for another VR or a point-to-point connection).

The origin node on the VR *sends* its bit-vector to its downstream neighbor. Upon receiving a bit-vector, a node performs a logical AND operation between the received and local vector (see Figure 10.b). The result is stored as the updated local bit-vector, and sent to the downstream node (see Figure 10.c).

When the origin receives a bit-vector from its neighbor (step 1), the first phase of the algorithm is completed. Note

Algorithm: *Feasible_Ring_Schedule* performed at node i

Input: VR_r ; origin node o_r ; p_r

Preprocessing: for each TF on this link
compute its RID (for example, o_r can be used as base)
construct BV_i to indicate the status of TFs

Begin

RECEIVE a bit-vector BV_{i-1}

1 **if** node $i = o_r$ then

1.1 **if** $BV_{i-1} == RBV$ then STOP

1.2 **else**

1.2.1 **if** number of 0's in $RBV < p_r$ then
print INFEASIBLE

1.2.2 **else**

$RBV = \text{MinMaxLinkSchedule}(VR, RBV, p_r)$
SEND RBV to downstream node $i + 1$

2 **else**

2.1 **if** $BV_{i-1} == RBV$ then $BV_i = RBV$

2.2 **else** $BV_i = BV_{i-1} \text{ AND } BV_i$

2.3 SEND BV_i to downstream node $i + 1$

End

Figure 11: Ring Scheduling Algorithm to Ensure Ring Feasibility.

that at this point each node knows what all the upstream nodes (from the origin to its upstream neighbor) know about available TFs. Thus, when the origin receives the bit-vector from its neighbor, it learns the ringwise availability of TFs (global information).

If the number of zeros in the received bit vector BV is less than p_r , the origin concludes (step 1.2.1) that the VR is *not feasible*. If BV contains at least p_r zeroes, then any choice of p_r TFs will result in a *feasible scheduling*. The origin calls the procedure $\text{MinMaxLinkSchedule}(VR, BV, p_r)$ (see algorithm in Figure 13), to choose p_r TFs for optimal scheduling (step 1.2.2). The bit-vector returned by this procedure indicates the TFs reserved for this multicast group. This bit-vector becomes the bit-vector of the reservation for this ring (RBV). The origin sends RBV to its downstream node and it is propagated on the ring to inform all the nodes. When the origin receives RBV (step 1.1), the algorithm terminates and each node on the ring has the same RID for the newly reserved TFs.

3.2 Optimal Selection of TFs

In this section we consider the problem of optimally choosing p_r TFs for a VR. First we introduce an objective function then an algorithm to construct a solution.

3.2.1 Objective Function to Minimize Jitter

Delay jitter determines buffer size at nodes and performance of real-time protocols. The larger the jitter the larger the buffers needed to compensate for it. A uniform distribution of TFs in the TC is beneficial in reducing the end-to-end delay, its variation (jitter), and the size of buffers in the buffering node. Ideally, if reserved TFs are uniformly distributed, the buffering node receives the packets at a (almost) regular

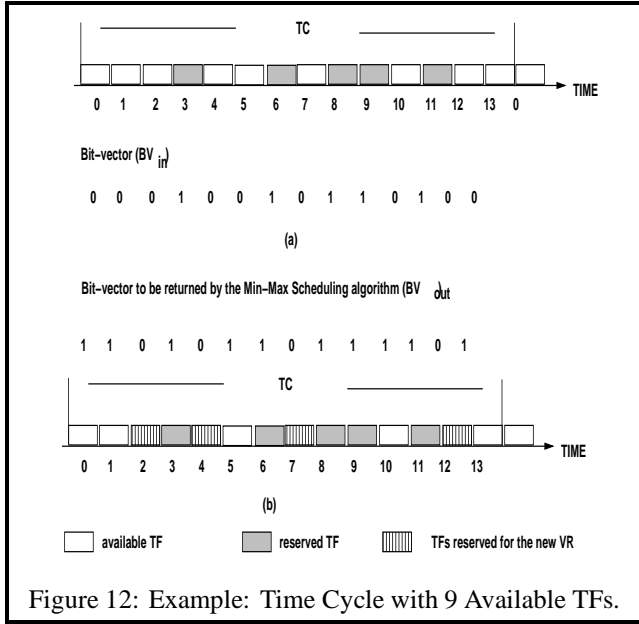


Figure 12: Example: Time Cycle with 9 Available TFs.

pace and forwards them at almost the same regular pace. As a consequence

1. packets are buffered shortly in the buffering node. Thus, the buffer must not be very large and the delay experienced by packets is small;
2. subsequent packets experience similar delays thus yielding small jitter.

Thus, one objective function for optimal schedule is to ensure equal *time distance* among TFs allocated in a TC. Let FS denote the set of all feasible schedules; we consider as optimal schedule σ the one which satisfies

$$\min_{\sigma \in FS} \left\{ \max_{1 \leq i \leq p_r} \left| [(t(\sigma)_i \bmod p_r + 1 - t(\sigma)_i + k) \bmod k] - \frac{k}{f} \right| \right\} \quad (3)$$

where $t(\sigma)_i$ denotes the i^{th} TF reserved for the multicast group according to the schedule σ .

Note that if all the TFs were available, the best spacing (timewise) between the p_r TFs would be: $\lfloor \frac{k}{f} \rfloor$. Thus, equation (3) compares the maximum difference between the spacing imposed by each feasible schedule to the best possible spacing and chooses the schedule which minimizes this difference.

Complexity of Link Scheduling

The size of the feasible scheduling set FS is determined by two parameters: p_r and $k - f$: enumeration can be obtained by $C(k - f, p_r)$ (i.e., number of ways to pick a subset of p_r out of $k - f$ distinct objects). Since the enumeration involves factorial operation, if p_r and $k - f$ are close to each other, it is possible to choose the best allocation by enumeration. However, for large factorial computations FS can be too large to enumerate; we propose a min-max algorithm to solve this problem with acceptable complexity.

Algorithm: *MinMaxLinkSchedule*

Input: VR_r ; BV_{in} ; p_r

Output: BV_{out} with exactly p_r zeros to indicate which TFs will be used by VR_r

Preprocessing: construct a weighted line-graph such that nodes represent available TFs
a link connects two consequent TFs
the weight (cost) on a link is the difference between RIDs of the end point TFs
mark the line-graph as “active” segment

Begin

1 Set all the bits to 1 in BV_{out}

Begin FOR $i \leq p_r$ **DO**

- 2 Choose the median of the active segment and reset the bit in BV_{out} identified by the median's index
- 3 Compute the total cost of each “active” segment
- 4 Store the costs and end points of the segments in list
- 5 Choose the segment with maximum total cost from this list and make this segment “active”

End

6 **Return**(BV_{out})

End

Figure 13: A TF Choice Algorithm to Minimize Jitter in a VR

3.2.2 Min-Max Algorithm

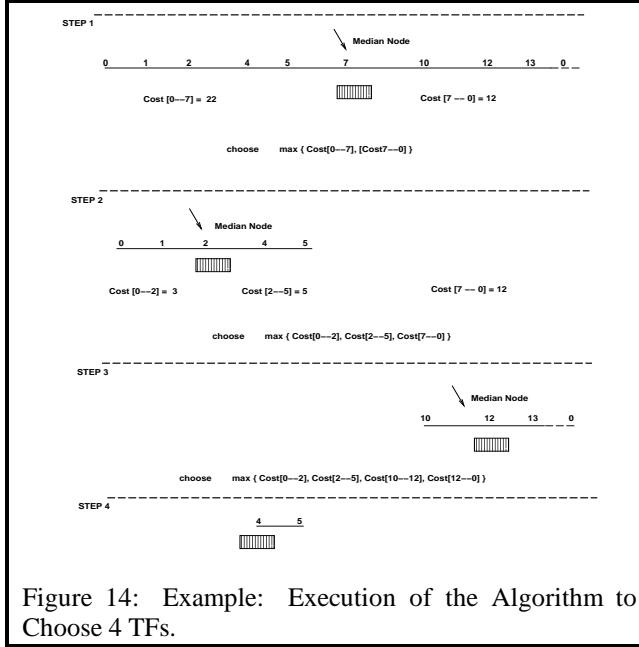
The input of the algorithm is a bit vector BV_{in} (Figure 12.a) which represents at least $k - f \geq p_r$ available TFs for a VR. The algorithm returns a bit vector BV_{out} (Figure 12.b) that indicates which TFs should be reserved for the multicast group.

The algorithm shown in Figure 13 builds a Line-Graph (LG) from the information carried by BV_{in} . Define a *median* of a LG as a node to which maximum *distance* from any other node in the graph is minimum. The *distance* between two nodes is defined as the sum of the link *costs* on the path between them as explained in Figure 13.

The algorithm takes p_r steps and at each step finds a median of an updated (reduced) graph as shown in Figure 14. Each median node corresponds to a TF selected by the algorithm for optimal scheduling. The location of median nodes in BV_{out} is identified by zero while all the other bits are set to 1. Each time a median node is chosen the LG is divided into two connected components (called segments). A new median is chosen from the segment with maximum total cost computed during the selection of the previous median node (see Figure 14). The selected segment is called the *active segment*; initially the entire LG is considered a single active segment.

3.2.3 Objective Function for Load Balancing

In this Section we relax the assumption that a TF is used to transmit a single packet. In this case different amounts of bits can be reserved during TFs. Balancing the load over the TFs allows a larger number of VRs to be allocated on the network since it decreases the blocking probability. Thus, we



introduce a new objective function which aims at balancing load when scheduling TFs for a VR. This objective function minimizes the differences among the amounts of non reserved bits in each of the k TFs of a TC. The optimal schedule is the one which satisfies

$$\min_{\sigma \in FS} \left\{ \sum_{i=1}^{p(\sigma)_r} \left(\beta(\sigma)_i - \overline{\beta(\sigma)} \right)^2 \right\} \quad (4)$$

where $\beta(\sigma)_i$ is the number of non reserved bits during the i^{th} TF, after the schedule σ has been applied. The variable $\overline{\beta(\sigma)}$ is the average amount of non reserved bits over the whole TC. The two objective functions (3) and (4) can be in conflict and thus a tradeoff between them must be found.

4 Adaptive Bandwidth Sharing

The adaptive bandwidth sharing concerns a single VR and does not affect any other VR embedded in the network. Active stations actually share the allocated bandwidth by sending their own packets during some of these TFs. Both passive stations and transit nodes just forward packets during the reserved TFs. The set of active stations changes dynamically as, for example, in a video-conference call the set of active speakers changes. The stations participating in a multicast group run a signaling protocol that allows each of them to know the actual set of active stations^{††}.

As the number a_r of active stations changes, the share of bandwidth available to each station changes and they have to (1) adapt their transmission rate to it (Section 4.1) and (2) agree on which TF each station is using (Section 4.2).

^{††}This signaling protocol is out of the scope of this paper; some possible candidate protocols are available for the purpose.

4.1 Rate Adaptation

Through one of the algorithms described in Section 4.2, each active station knows the amount of bandwidth it can use for transmission. When the bandwidth available to an active station changes, the application is notified in order to change the rate at which data is produced.

Some applications require some time in order to make this change. For example, in a video-conferencing with compression data is produced through a process starting with capturing video frames, going through a number of processing steps, and ending with quantization and run-length encoding of the video information. The compressed data is buffered and retrieved from the buffer at a constant bit rate. Changing the amount of data produced by this process implies tuning the parameters that regulate some of these processing steps. The earlier the involved processing step, the longer the time for having the change take effect at the network interface.

When a station becomes passive, it notifies other stations and immediately stops transmitting. Each active station realizes that its share of the bandwidth is larger and notifies the application to increase the data generation rate. Even though the increment of the bandwidth share is immediately available to the active nodes, it takes some time before the change in the application generation rate has effect and is propagated to the network interface; i.e., the transmission rate of active nodes does not increase instantly.

When a station becomes active, it notifies other stations. However, it does not start transmitting immediately. All the previously active stations immediately notify their application to lower the data generation rate. Since it takes some time before the actual transmission rate of all the active stations decreases, the new active station waits for its bandwidth share to be available before starting transmitting. We assume that newly active stations wait a fixed time (in the order of a second) before transmitting. This is needed in order to assure that no packet is dropped in either the network interface of active stations or inside the network, due to unavailability of communications resources.

This waiting time is acceptable as the time scale in which stations change their activity status is usually long. For example, in a video-conferencing application a station becoming active implies that somebody is asking the right to speak. Thus, waiting should not be of concern since the new speaker will usually have to wait before getting the permission to speak.

4.2 Bandwidth Sharing

Given the set of the p_r TFs reserved to a multicast group r , the bandwidth sharing problem consists in having each active station using a subset of them to send its own packets^{***}. The subset used by each active station should be built satisfying the following criteria:

^{***}Whenever dealing with TFs in this Section, we mean TFs reserved to the VR.

1. each node gets a fair share of the bandwidth MB_r allocated to the multicast group;
2. the end-to-end jitter experienced by packets is minimum;
3. the unused bandwidth (i.e., the TFs during which no multicast data is transmitted) is minimized*.

4.2.1 Round Robin Allocation

The allocated bandwidth can be shared fairly among the active stations by having each active station using each TF in a round robin fashion. Active stations transmit their packets in any *free* or unused TF. When a station receives a packet originated by itself, it must not use the next TF to send another packet of its own. This allows downstream stations to use the TF, in a round robin manner, for their own transmissions. The fraction of unused bandwidth due to this phenomenon is:

$$\frac{1}{a_r + 1}$$

This unused bandwidth is analogous to the unused bandwidth due to the token walk time in token ring.

4.2.2 Fixed Allocation

A fixed allocation scheme can be used to reduce the amount of bandwidth unused due to the bandwidth sharing procedure. According to the fixed allocation scheme, active stations come to an agreement on which TFs each of them is using[†]. When a station receives back a packet sent by itself, it can use the next TF for sending its data. Figure 15 shows an example in which 4 TFs (RID 1, 4, 8, 12) have been reserved for the VR depicted in Figure 5 and nodes B and C are active. Each active node is using two TFs to transmit its data.

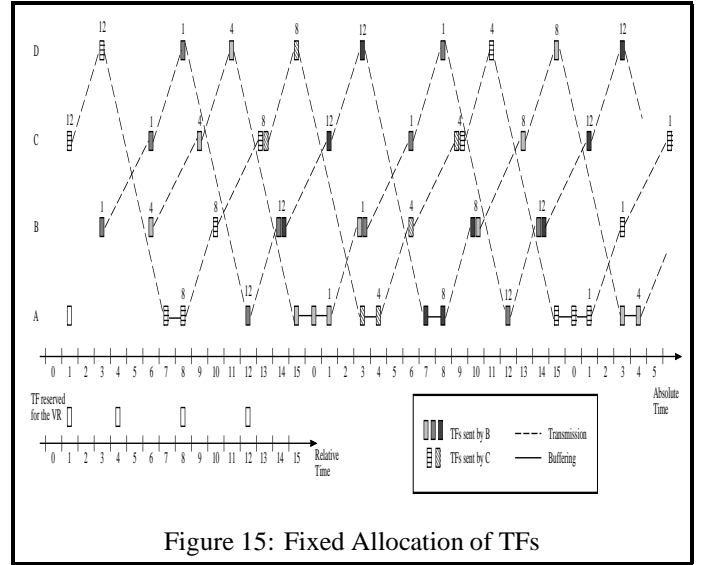
Note that the RIDs of the two TFs each node is using change over time because the buffering node A buffers and forwards packets in a different TF. Moreover, some bandwidth is still unused because in each moment only 3 TFs out of the 4 reserved are used. This is due to the behavior of the buffering node, the distribution of the TFs in the TC, and the relationship between the lengths of VR and TC. This situation can be accepted because with this behavior the buffering node delays packets for a short time as shown in the example of Figure 15.

In the example shown in Figure 15 the phenomenon is particularly harmful as the unused TF is always one of the 2 assigned to station C which is getting only half of the fair bandwidth share.

This unfairness is much smaller as larger is the number of TFs assigned to each station. This unfairness can be avoided by applying more complex bandwidth sharing schemes or buffering policies in the buffering node. When dealing with

* Even though unused bandwidth can be used to transmit best effort traffic, the VR should try and use as much as possible of the bandwidth allocated to it.

[†] This can be done through some signaling mechanism that is out of the scope of this paper.



unused bandwidth in the following, the contribution given by the foregoing phenomenon is not taken into account.

The fixed allocation of TFs to active stations is possible only if $a_r \leq p_r$, otherwise a different scheme (e.g., round robin allocation) must be applied.

Nevertheless, one of the main fields of application for real-time group multicasting is large scale video-conferencing, where it is possible to have a large number of listeners (i.e., large p_r) and a small number of active speakers (i.e., small a_r). Only the active speakers can transmit to the group.

When $p_r \bmod a_r = 0$, fixed allocation of TFs lets each active station use $\frac{p_r}{a_r}$ TFs per TC, i.e., active stations share the bandwidth fairly and do not leave unused any fraction of it. The jitter introduced depends on the distribution of TFs in the TC and the allocation of TFs to active stations.

When $p_r \bmod a_r \neq 0$, p_r can be written as

$$p_r = I \cdot a_r + N, \text{ where } I = \lfloor \frac{p_r}{a_r} \rfloor \text{ and } N = p_r \bmod a_r$$

Each active node uses the fixed allocation criteria on I TFs; for the remaining N TFs several approaches can be followed.

1. The N TFs are allocated to the active nodes in a round robin fashion. As a result, with this approach some bandwidth is not used and a large jitter is introduced because usually a node sends I packets per TC and seldom it sends $I + 1$ packets. This exceptional packet received between two regular packets can sharply reduces the packet inter-arrival time, thus yielding a large variation of the end-to-end delay (increase of the delay jitter).
2. N of the a_r active nodes fixedly use $I + 1$ TFs for sending their packets. The bandwidth is not shared fairly, but all the allocated bandwidth is used and the bandwidth sharing scheme does not contribute to the jitter.
3. The N TFs are left unused; the multicast group does not use some of the bandwidth allocated to it, but the

bandwidth allocation scheme does not contribute to jitter and it is very simple.

5 Conclusions

The combination of virtual rings and timing information is interesting and intriguing. In this work we have shown how useful this combination can be for *adaptive real-time multicast*, while in a previous work we have shown that it is also useful for *reliable multicast* for data processing applications [4].

The proposed approach to adaptive real-time multicast is characterized by having a *deterministic behavior*. This is important for real-time applications in two ways: (1) it minimizes the *bound* on the delay and delay variation across each virtual ring, and (2) it ensures no loss due to congestion. The no loss property is important for real-time video applications using compression. In compression a loss of a single packet can result in the inability to properly display multiple consecutive video frames, which can annoy the viewer.

Multicast communication over embedded virtual rings has a drawback that the delay is increased with the ring physical size and number of participants. Approaches based on trees do not suffer from this problem. However, tree based multicast suffers from two basic problems: (i) loss due to congestion, and (ii) queuing delay bound that can be very large. It is interesting to investigate ways to exploit some of the tree properties while maintaining the virtual rings virtues: (i) fair sharing of the bandwidth among the active nodes in the multicast group, (ii) high and stable utilization of the bandwidth, and (iii) deterministic performance with (very) small packet loss and (very) small queueing delay bound.

There are several open issues which require more research. The scheduling and allocation of time frames to the various multicast groups (virtual rings) can be improved by further investigating how to reduce the blocking probability. Blocking is defined as the impossibility of allocating enough bandwidth for a new multicast group while some capacity is still available, but not in the proper time frames. Another issue concerns the unused capacity assigned to a virtual ring due to the bandwidth adaptation schemes proposed in this work.

References

- [1] P. H. Dana. Global positioning system overview. <http://www.utexas.edu/depts/grg/gcraft/notes/gps/gps.html>.
- [2] C-S. Li, Y. Ofek, A. Segall, and K. Sohraby. Pseudo-isochronous cell switching in atm networks. *Proc. IEEE INFOCOM'94*, pages 428–437, 1994.
- [3] C-S Li, Y. Ofek, and M. Yung. “Time-driven priority” flow control for real-time heterogeneous internetworking. *IEEE INFOCOM'96*, 1996.
- [4] Y. Ofek and B. Yener. Reliable concurrent multicast from bursty sources. *To Appear in IEEE JSAC (preliminary version is in Proc. IEEE INFOCOM'96)*, 1996.