# Rate Based Congestion Control for Multicast ABR Traffic

D. Cavendish[1], S. Mascolo[2], and M. Gerla[1] *

1.Department of Computer Science
University of California, Los Angeles CA 90024 - USA
{dirceu,gerla}@cs.ucla.edu

2.Dipartimento di Elettrotecnica ed Elettronica
Politecnico di Bari, 70125 Bari - Italy
mascolo@poliba.it

*Abstract* - In this paper we present an end to end rate control algorithm for multicast (point to multipoint) ABR service in ATM networks. The algorithm is the multicast extension of the SP-EPRCA unicast (point-to-point) congestion control algorithm proposed in [1]. The goal is to control the multicast input rate in order to achieve high bandwidth utilization without overflowing any queue along the multicast tree. The straightforward approach would be to consider a multicast connection as a superposition of unicast connections, and to compute the multicast input rate as the minimum among all unicast input rates. The main problem is the congestion caused by feedback traffic from the multiple destinations. To solve this problem, we show how feedback information cells can be intelligently "merged" at each multicast tree "fork" without losing essential information needed by the control algorithm placed at the source. The algorithm proposed inherits all properties from the unicast SP-EPRCA.

## I  Introduction

Congestion control is known to be one of the main challenges of wide area high speed networks, such as ATM, which are characterized by high bandwidth-delay products. Congestion is either avoided by preallocating bandwidth (CBR and VBR services), or is controlled by regulating input traffic rate at the source (ABR and UBR). In this paper, we address rate based congestion control approach only.

As recommended by ATM Forum, the rate based congestion control framework for ABR/VBR traffic follows a closed-loop control scheme, where the input rates are adjusted based on current network congestion levels [5]. The congestion information is collected by RM (resource management) cells, and transported all the way back to the source, to perform rate adjustement. The challenge is, however, that the large delay inside the feedback loop due to the large bandwidth-delay product affects stability. Fitting into this closed-loop framework, many algorithms have been recently proposed. In particular, we have proposed the Smith Predictor Enhanced Proportional Rate Control Algorithm (SP-EPRCA) [1], which shows how to take into account the feedback loop delays when input rates are computed. The main advantage

of SP-EPRCA over other rate based schemes is to strictly prevent cell loss when adequate buffers are allocated to each connection at each intermediate node.

Cell loss prevention is very desirable in reliable multicast applications which require detection and recovery of lost data (e.g., distributed simulation). Note that while unicast is generally protected by TCP, multicast runs under UDP (which does not provide retransmissions of lost packets). Thus, in ATM, the loss of a multicast cell must be recovered at the application layer, with higher latency and processing overhead than at the transport layer [2]. If the cell loss is drastically reduced or better yet, eliminated (by the SP-EPRCA control algorithm, for example), then the payoff performance vs cost is substantial.

For multicast ABR services, it is important to investigate the additional difficulty introduced by multicast rate controlling. Very few researchers have addressed the problem so far. A significant contribution is by *Kay et al.* [4], which shows how it is possible to extend a max-min fair unicast explicit rate indication algorithm to multicast traffic.

This paper extends the SP-EPRCA algorithm to multicast traffic. The paper is organized as follows: We first consider a multicast connection as a collection of independently rate controlled unicast connections and define an ideal multicast rate control algorithm. Then we propose a way of "merging" RM cells so that the feedback traffic be reduced back to the forward control traffic levels. The cell loss free and min-max fairness properties come primarily from the unicast SP-EPRCA algorithm. Simulations are used to substantiate our claims.

## II  Ideal Multicast Rate Control Algorithm

In this section we first present the unicast SP-EPRCA algorithm in which the multicast rate control algorithm is based upon. Then we define the ideal multicast rate control algorithm(IMRCA), as well as describe its major implementation difficulty.

### II.1  Unicast SP-EPRCA Algorithm

The continuous time model of the controlled system for unicast connections is modelled as shown in Fig. 1. There we represent the bottleneck queue of a VC by an integrator

(fluid flow model), the Smith Predictor controller by $K^*$, and the delays $T_{fw}$ and $T_{fb}$ inside the control loop. The variables of interest are:

$u(t)$ - Source input rate
$x(t)$ - Bottleneck buffer level
$d(t)$ - Bottleneck service rate
$X^o$ - Target filling level of switch queues
$K$ - Gain of the proportional controller
$T_{fw}$ - Feedforward delay
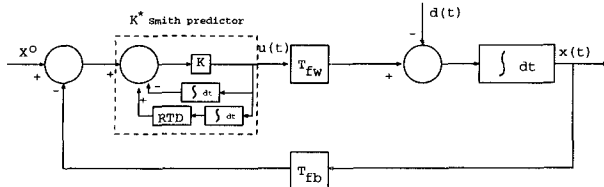$T_{fb}$ - Feedback delay



Figure 1: Queue dynamic model using a proportional controller plus a Smith Predictor

The control equation is given by:

$$u(t) = K[X^o - x(t - T_{fb}) - \int_{t-RTD}^{t} u(t')dt']\qquad(1)$$

This equation can be interpreted as a proportional control of the queue level $N(t) = x(t) - \int_{t-RTD}^{t} u(t')dt'$, i.e. the in pipe cells can be seen as being stored at the bottleneck queue. *Mascolo at al*[1] presents a discrete time version of this continuous time control equation, in order to provide the effective equation used by the controller. Requirements about sampling frequency arises, which lead to the creation of the virtual feedback procedure to ensure loss free property in case of missing feedback information. The reader is referred to [1] for further information.

## II.2 IMRCA Definition

**Definition 1** *The ideal multicast rate control algorithm (IMRCA) sets the rate of a multicast connection $(S, D^1, D^2, \cdots, D^M)$ to the minimum of the rates computed by M unicast SP-EPRCA controllers executed independently for each pair $(S, D^i)$, where each algorithm is fed by its corresponding RM cell stream. Before the multicast rate updating, the source must receive one RM cell from each multicast destination. In case any of RM cells is missing for a time interval longer than $\triangle$, the algorithm performs virtual feedback [1] across individual unicast SP-EPRCAs in order to update the multicast input rate as the minimum among them.*

In other words, each unicast controller placed at source site computes its rate $u_i$ according to equation (1). Since, by definition of multicast, the input rate is unique across the entire multicast tree, each unicast controller differs from the others by the round trip delay (RTD) only. When a unicast rate $(u_j)$ is computed, say at time $t$, the multicast rate is then updated by:

$$u_{min}(t) = \min_{1 \leq j \leq M}(u_j)\qquad(2)$$

The straightforward implementation of the ideal multicast algorithm is as follows. Data and RM cells are generated at source site (at a proportion of NRM data per RM cell [1]), and travel downstream along the multicast tree, getting replicated at fork nodes (see Fig. 2). At each fork of the multicast tree, RM cells are stamped with their respective destinations IDs when replicated so that each RM stream can be identified at the source site. In the backward direction, RM cells are sent upstream to the respective source. In other words, $M$ SP-EPRCA algorithms are running independently, each one fed by its own RM cell stream.



Figure 2: Splitting at multicast tree fork

Although this algorithm is clearly ideal (IMRCA), it is not feaseble, because the control traffic that is generated downstream the tree must return upstream, most likely causing congestion. As an additional drawback, the computation burden at the source grows with the number of destinations. In the next section we show how to solve this problem.

## III Merging RM cells

The previous section has exposed the infeasibility of sending multiple control streams back to the source for rate control. Clearly the main issue here is how to "merge" RM cells at fork nodes so that useful feedback information is *not* lost on their way back to the source. In essence, merging should be the inverse process of splitting. If a forward stream has been splitted in $K$ streams at fork $F$, then the $K$ control feedback streams are merged into one stream at $F$. Therefore, the merging process requires that the fork $F$ awaits for one RM cell from each of its "children" before deciding what to send back in a single RM cell. Figure 3 shows the RM merging at a tree fork.



Figure 3: Merging control traffic at multicast tree fork

Suppose all $K$ RM cells were able to reach the source, so that the source runs independent unicast SP-EPRCA controllers, one per source/destination path. Let $U$ be the vector of rates $U = (u_1, u_2, \cdots, u_M)$ computed. It is clear that, using IMRCA, only one SP-EPRCA controller will dictate the next rate to be enforced. Therefore, the idea is to detect, at each fork point, which RM cell information will result in the lowest rate, and pass only that value upstream.

### III.1 How to choose the right RM cell information
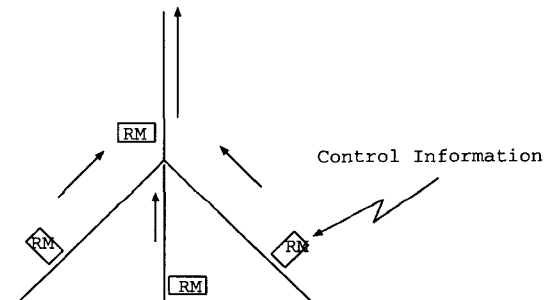
Considering a generic fork $F$, let $td_f$ be the transmission delay from $F$ to the source. Assume that at time $t$ the source gets one $RM_i$ cell from each $D_i$. Then the unicast rate $u_i$ will be computed as [1]:

$$u_i = K[X^o - x_i - \int_{t-RTD}^{t} u(\tau)d\tau] = K[X^o - N_i]^1 \quad (3)$$

where $N_i = x_i + \int_{t-RTD}^{t} u(\tau)d\tau$.

A physical interpretation of equation (3) is that the new rate is proportional to the available space at the bottleneck buffer ( which has size $X^o$). The available space is $X^o$ minus the queue level information received in the RM cell, minus all the in flight cells in the feedback loop, i. e. $N_i$, which might get stored in this buffer.

The minimum rate is, thus, given by the $j$th controller (corresponding to destination $D_j$) such that:

$$x_j + \int_{t-RTD_j}^{t} u(\tau)d\tau = \max_{1 \leq i \leq M}[x_i + \int_{t-RTD_i}^{t} u(\tau)d\tau] \quad (4)$$

It remains to be shown how fork $F$ can decide which of its children will provide the feedback information that satisfies Eq. (4) at the time $t' = t - td_f$

We can write $N_i$ as follows:

$$N_i = x_i + \int_{t-RTD_i}^{t-td_f} u(\tau)d\tau + \int_{t-td_f}^{t} u(\tau)d\tau = x_i + c_i + D \quad (5)$$

Since $D$ is constant for all $i$, the $j$ that maximizes $N_i$ is given by the $j$ that maximizes $(x_i + c_i)$.

Therefore, at time $t'$, fork $F$ selects the $j$th feedback value such that:

$$x_j + c_j = \max_i(x_i + c_i) \quad (6)$$

This computation is performed sequentially from the destinations (multicast tree leaves) up to the source.



Figure 4: Feedback Information Selection

### III.2 The SP Multicast Rate Control Algorithm (SP-MRCA)

The fork nodes in the multicast tree perform the merging scheme described previously, as long as they receive RM values from all their children in a timely fashion. Since at multicast setup time no information is present in the tree, the algorithm does not receive feedback information until $RTD_{max}$ has elapsed, where $RTD_{max} = \max_{1 \leq i \leq M}(RTD_i)$ is the maximum round trip delay among all $D_i$ belonging to the multicast connection. Therefore, we need to specify the algorithm behavior for two cases:

i) $t < RTD_{max}$: Since the algorithm does not receive feedback cells, it uses the virtual feedback of the unicast SP-EPRCA algorithm with $RTD_{max}{}^2$.

ii) $t \geq RTD_{max}$:

1- Upon receiving RM cells, the algorithm updates the rate using the unicast SP-EPRCA controller corresponding to the destination $D_i$ from which the cell comes from.

2- Upon expiration of the feedback interarrival time $\triangle$, the source uses the last RM cell received and the corresponding unicast SP-EPRCA algorithm to perform a virtual feedback[1].

**Theorem 1** *The SP-MRCA algorithm is equivalent to the IMRCA algorithm.*

We have to show that the two algorithms compute always the same rate. Since the proof is rather intuitive, for sake of brevity we give only a sketch of it.

**Sketch of Proof 1** *We divide the proof in two cases:*

*i) $t < RTD_{max}$: During this interval, the connection with longest RTD does not receive feedback information, and*

---

[1]We use $x_i$ instead of $x_i(t - t_{fb})$ for sake of conciseness

[2]We assume a $RTD_{max}$ large enough such that virtual feedback is needed, since the other case is trivial (see [1])

*thus uses virtual feedback for its rate computation, say $u_v$. The rate $u_v$ results to be the minimum among all independent SP-EPRCA controllers, and hence is equal to the rate computed by the IMRCA.*

*ii)* $t \geq RTD_{max}$ :

    *1- Upon receiving RM cells, our merging scheme at the fork points sends upstreams the RM cell corresponding to the minimum rate. Thus the theorem follows in this case.*

    *2- Upon expiration of the feedback interarrival time $\triangle$, the IMRCA uses virtual feedback, that is, it increases each $x_i + c_i$ by the same amount (since the multicast rate is unique). Therefore, the minimum rate decreases, but it results to be computed by the same $j$ controller as in the previous iteration, as it is in the SP-MRCA.*

## III.3  A pratical way for forks to filter out excessive RM cells

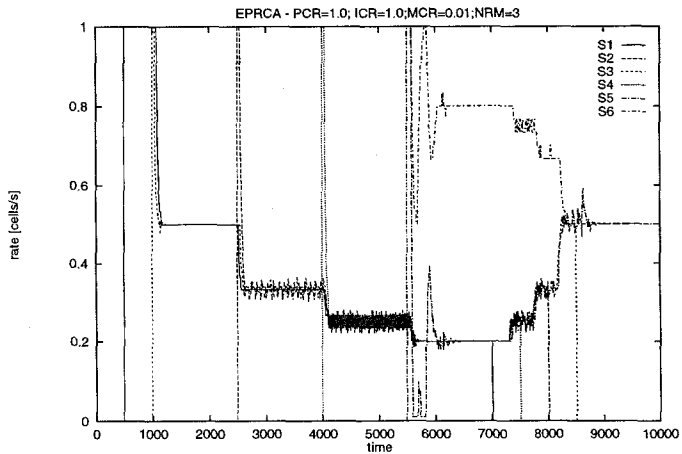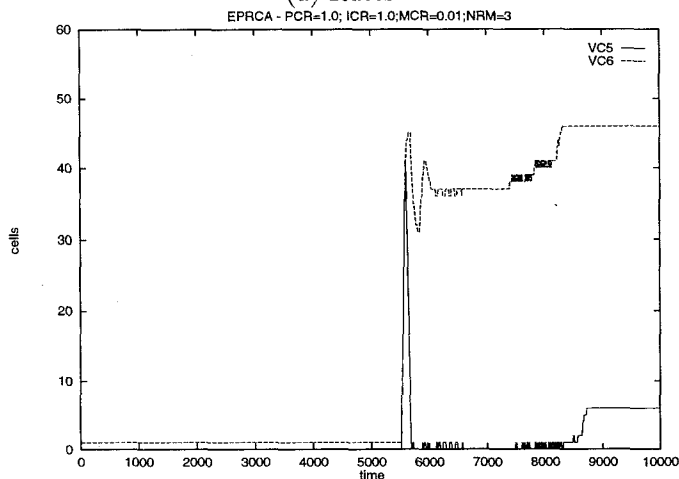In merging feedback information, cells must be filtered out at fork points. In the SP-MRCA algorithm, a fork must have feedback information about all its children before relaying any information upstream the tree. In principle, feedback information from shorter paths must wait for that coming from the longest path. This suggests the following scheme for RM cell filtering and information reporting back to the source: Define the longest path $(F, D_j)$ among all $D_i$ from a fork node $F$ to be a kind of "conveyor belt" of feedback information. That is, RM cells circulating in this path are used to relay feedback information upstream after merging. Namely, all cells coming from a branch other than the conveyor belt are filtered out after having their feedback information recorded. The merged information then awaits for the next RM cell coming from the conveyor belt to be carried back to the source. Each fork has its own conveyor belt, which is determined by its distance to its destinations (sub tree leaves rooted at $F$). The scheme is shown in Figure 5.



Figure 5: Recursive filtering of RM cells

## IV  Simulation Results

The network topology used in a simulation of the SP-MRCA is shown in Fig. 6. Links have uniform speeds, normalized to 1 cell per unit of time [cells/s]. Link labels represent the bandwidth-delay product. There are five unicast connections: $(S1, D1),(S2, D2),(S3, D3), (S4, D4),(S6, D6)$. There is also one multicast connection $(S5, D5, D7, D8)$, spanning across multiple bottlenecks.



Figure 6: Network Topology

Connections are set and released according to Table 1. This time epochs were chosen so that bottlenecks shift in time from link to link across the network. We assume infinite backlog at each source. Buffer capacity is 60 cells.

Table 1: VC Connections Activity

| Connection # | 1 | 2 | 3 | 4 | 5m | 6 |
|---|---|---|---|---|---|---|
| Start Time | 500 | 2500 | 1000 | 4000 | 5500 | 0 |
| End Time | 7000 | 8000 | 8500 | 7500 | 10000 | 10000 |
| RTD (cells) | 0 | 20 | 40 | 40 | 80(max) | 0 |

Figure 7(a) shows the behavior of the six input rates, corresponding to connections $S1 - S6$, at source nodes. We assume an initial cell rate of 1 [cell/s], although the results are independent of this value. After each traffic pattern change ( start/end of a connection), each rate rapidly settles to a new fair stationary value. Sources $S1$, $S2$, $S3$, $S4$, and $S5$ (multicast) are constrained by bottleneck #2, and thus get the same bandwidth share each, while source $S6$ is constrained by bottleneck #1, and therefore takes all the bandwidth left by the multicast source $S5$. We can see, also, that due to traffic variations during simulation time, the multicast bottleneck moves along the tree (i.e., from bottleneck #2 to bottleneck #1, see Figure 6). Figure 7(b) shows the dynamic behavior of the two queues at bottleneck #1, corresponding to $VC5, VC6$, while Fig. 7(c) shows the dynamic behavior of the five queues at bottleneck #2, corresponding to $VC1 - VC5$ bottleneck queues (since connection #6 does not cross this bottleneck). No queue overflow occurs, by virtue of the SP-EPRCA overflow prevention property. Figure 7(c), for time $> 8500$, shows all queues emptying due to traffic termination, which again means that the bottleneck queue for multicast source $S5$ has moved from bottleneck #2 to bottleneck #1 (Fig. 7(b)). The overall performance is consistent with the ideal multicast rate control algorithm (IMRCA), as expected.
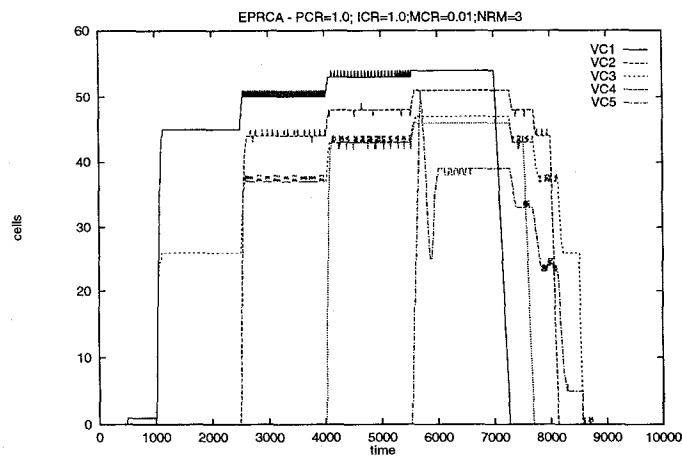
(a) Rates



(b) VC Queues at Bottleneck #1



(c) VC Queues at Bottleneck #2

Figure 7: Multicast Rate Control Algorithm

## V  Conclusion

We have proposed an efficient multicast rate control algorithm (SP-MRCA), which extends the unicast SP-EPRCA algorithm to multicast traffic. Key issues, such as the reduc-

tion of excessive feedback information, were addressed. Simulation results were used to exemplify the algorithm performance with respect to dynamic tracking of bottlenecks, cell loss free property and fairness. The trade-offs between buffer sizes and throughput are the same as in the SP-EPRCA, and thus were not addressed here. Although we have assumed individual queues per VCs at each switch output port, work is in progress to provide single buffer per output port, as reported in [3] for unicast connections. As for the implementation of the Smith Predictor controller at the source, this can be placed either in the Host or in the edge switch. The latter implementation is preferred since it complies with the ATM Forum UNI specifications. In this case, the entire SP-MRCA is transparent to the end user; the edge switch simply dictates the allowed rate to the Host.

It is worth saying that the algorithm can be easily extended to support dynamic join/leave of receivers in the multicast group. Namely, if the multicast VC tree can be dynamically expanded/contracted, then the SP-MRCA algorithm will automatically adjust to the new configuration. Likewise, the SP-MRCA can also be extended to the many-to-many multicast case.

Strict buffer overflow protection via SP-MRCA may require large buffers for WAN type propagation delays in order to have full bandwidth utilization. To reduce implementation cost, without giving up bandwidth utilization, smaller buffers than the strictly required for large delays may be used, thus incurring some cell loss. In this case, Host retransmissions using schemes such as the one reported in [2] must be used. Work is in progress to determine the trade-offs between buffer costs and Host retransmission overhead and latency.

## References

[1] S. Mascolo, D. Cavendish, M. Gerla,"ATM Rate Based Congestion Control Using a Smith Predictor: an EPRCA Implementation", *Proceedings of INFO-COM'96*, vol.2, pp.569-576.

[2] S. Floyd, V. Jacobson, S. McCanne, C. G. Liu, and L. Zhang, "A Reliable Multicast Framework for Lightweight Sessions and Application Level Framing," *Proceedings of SIGCOMM'95*,pp.342-356.

[3] D. Cavendish, M. Gerla, and S. Mascolo "ATM Rate Based Congestion Control Using a Smith Predictor: Implementation Issues", *Proceedings of First Workshop on ATM Trafic Management, IFIP*,pp. 289-296.

[4] K. Y. Siu and H. Y. Tzeng, "Congestion Control for Multicast Service in ATM Networks", *Proceedings of GLOBECOM'95*,pp. 310-314.

[5] A. W. Barnhart, "Baseline Performance Using PRCA Rate-Control," *ATM Forum/94-0597*, July 1994.