

Rate-Controlled Static-Priority Queueing*

Hui Zhang

Domenico Ferrari

hzhang, ferrari@tenet.Berkeley.EDU

**Computer Science Division
University of California at Berkeley
Berkeley, CA 94720**

TR-92-003

February 1992

Abstract

We propose a new service discipline, called the Rate-Controlled Static-Priority (RCSP) queueing discipline, that can provide throughput, delay, delay jitter, and loss free guarantees in a connection-oriented packet-switching network. Previously proposed solutions are based on either a time-framing strategy, or a sorted priority queue mechanism. Time-framing schemes suffer from the dependencies that they introduce between the queueing delay and the granularity of bandwidth allocation; sorted priority queue may be difficult to implement. The proposed RCSP queueing discipline avoids both time-framing and sorted priority queue; it achieves flexibility in the allocation of delay and bandwidth, as well as simplicity of implementation. The key idea is to separate rate-control and delay-control functions in the design of the server. Applying this separation of functions will result in a class of service disciplines, of which RCSP is an instance.

* This research was supported by the National Science Foundation and the Defense Advanced Research Projects Agency (DARPA) under Cooperative Agreement NCR-8919038 with the Corporation for National Research Initiatives, by AT&T Bell Laboratories, Hitachi, Ltd., Hitachi America, Ltd., Pacific Bell, the University of California under a MICRO grant, and the International Computer Science Institute. The views and conclusions contained in this document are those of the authors, and should not be interpreted as representing official policies, either expressed or implied, of the U.S. Government or any of the sponsoring organizations.

1 Introduction

Future high speed networks will have to support real-time communication services. A real-time communication service will allow clients to transport information with performance guarantees expressed in terms of delay, delay jitter, throughput and loss rate bounds. It has been argued that a connection-oriented architecture, with explicit resource allocation and connection admission control, is needed to offer such a real-time service [6]. However, in a packet-switching network, packets from different connections will interact with each other at each switch; without proper control, these interactions may adversely affect the network performance experienced by clients. Service disciplines of the switching nodes, which control the order in which packets are serviced, determine how packets from different connections interact with each other.

The First-Come-First-Served (FCFS) discipline, which is used in most conventional networks, may be capable of providing some performance guarantees, but only under a number of constraints and with a very inefficient use of the network's resources [5]. New service disciplines that aim to provide different qualities of service to different connections have been proposed. In a previous study [15], six rate-based service disciplines were compared: Delay Earliest-Due-Date (Delay-EDD) [6], Virtual Clock [16], Fair Queueing [2], Jitter Earliest-Due-Date (Jitter-EDD) [14], Stop-and-Go [7], and Hierarchical Round Robin (HRR) [9]. It has been shown that the three non-work-conserving disciplines¹, Jitter-EDD, Stop-and-Go and HRR, exhibit features that are suitable for offering guaranteed network performance: they all maintain some sort of *traffic smoothness* inside the network, and therefore the buffer space they need to prevent packet loss is uniformly distributed over the switches along the path of each connection; they can all provide guaranteed throughput and bounded-delay services; Jitter-EDD and Stop-and-Go can also provide bounded-delay-jitter service. However, each of them suffers from some disadvantages. In particular, both Stop-and-Go and HRR use a framing strategy, which introduces dependencies between delay and bandwidth allocation granularity, so that connections with both low delay and low throughput requirements can not be supported efficiently. Jitter-EDD needs to implement a sorted priority queue, which has an $O(\log N)$ insertion operation [11], where N is the number of packets in the queue. This may not be feasible for a high-speed implementation. Jitter-EDD also requires a complicated schedulability test at connection establishment time [6, 15], which is not needed in Stop-and-Go and HRR.

In this paper, we will present a new service discipline called *Rate-Controlled Static-Priority (RCSP) queueing* that overcomes these limitations. In particular, RCSP

- provides delay, delay jitter, throughput and loss-free guarantees;
- requires evenly distributed amounts of buffer space at the switches to prevent packet loss;
- introduces no coupling between delay and bandwidth allocation;
- is relatively easy to implement, since no sorted priority queue is needed; and
- requires relatively simple admission control tests.

In section 2, we review some of the issues that arise when designing service disciplines to provide performance guarantees. In section 3, we describe the algorithm of the proposed RCSP discipline and the corresponding admission control conditions. In section 4, the new discipline is compared with the previously proposed Jitter-EDD, Stop-and-Go, and HRR disciplines. We show that the key idea is to separate the functions of rate-control and delay-control; this separation of functions will result in a class of service disciplines, of which RCSP is an instance. Finally, we give conclusions and outline future work in section 5.

¹ In a work-conserving discipline, a server is never idle when there is a packet to transmit; in a non-work-conserving discipline, each packet is assigned, either explicitly or implicitly, an *eligibility time*; even when the server is idle, if no packets are eligible, none will be transmitted [15].

2 Background

We consider the paradigm proposed in [6] for providing guaranteed service to clients in a packet switching network: before communication starts, the client specifies its traffic characteristics and performance requirements to the network; the client's traffic and performance parameters are translated into local parameters, and a set of connection admission control conditions are tested at each switch; the new connection is accepted only if its admission would not cause the performance guarantees made to other connections be violated; during data transfers, each switch will service packets from different connections according to a service discipline; by ensuring that the local performance requirements are met at each switch, the end-to-end performance requirements can be satisfied. Notice that there are two levels of control in this paradigm: connection admission control at the connection level, and service discipline at the packet level. A complete solution needs to specify both the service discipline and the associated connection admission control conditions.

To design a service discipline that guarantees performance, two issues need to be addressed:

1. how to control the interactions between different connections in a single switch;
2. how to control the interactions between different switches.

	Bandwidth Allocation	Delay Allocation
Virtual Clock		
Fair Queueing		Sorted Priority Queue
Delay-EDD		
Jitter-EDD	Regulator	Sorted Priority Queue
Stop-and-Go		Multi-level Framing
Hierarchical Round Robin		

Table 1: Mechanisms to Allocate Delay and Bandwidth

The first issue relates to the design of a single server: how to allocate delay and bandwidth among different connections. Previously proposed solutions can be summarized as shown in Table 1.

As stated before, a multi-level framing strategy introduces dependencies between delay and bandwidth allocation. A sorted priority queue avoids this coupling, but may be difficult to implement in very high-speed switches; also, to provide a continuous spectrum of delay bounds, it needs a complicated schedulability test at connection establishment time.

The second issue concerns the interaction between different switches along a path. A switch can provide local performance guarantees to a connection only when the traffic on that connection is well-behaved. However, network load fluctuations at previous switches may distort the traffic pattern of a connection and cause an instantaneous higher rate at some switch even when the connection satisfies the client-specified rate constraint at the entrance to the network. If the maximum distortion of the traffic patterns can be bounded, buffer space can be allocated at each switch to accommodate the distortion. Delay-EDD adopts this approach. Unfortunately, distortions of traffic patterns due to network load fluctuations tend to accumulate, and this requires the amount of buffer space allocated to grow linearly along the path. Buffer space can be saved by controlling the distortion of traffic pattern at each switch. In Jitter-EDD, Stop-and-Go, and HRR, traffic patterns are partially or

completely reconstructed at each switch so as to offset the effects of network load fluctuation and interaction between switches.

A relevant issue here is that of jitter control. In the literature, there are different definitions of *jitter*. In [9], the term is used to capture the *burstiness* of the traffic, and is defined to be the maximum number of packets in a *jitter averaging interval*. In [3, 14], the term is used to capture the magnitude of the distortion to the traffic pattern caused by the network, and is defined to be the maximum difference between the delays experienced by any two packets on the same connection. In this paper, we call these two quantities *rate jitter* and *delay jitter*, respectively. Eliminating delay jitter at each switch along the path will *completely* reconstruct the traffic pattern at that switch; thus, if the traffic obeys the traffic specification at the entrance to the network, it will obey the specification throughout the network. Controlling the rate jitter at each switch along the path will only *partially* reconstruct the traffic pattern; the partially reconstructed traffic pattern will obey the same traffic specification as the input traffic, but some inter-packet spacing information is lost. Notice that, if the traffic at the entrance to the network has a small rate jitter, controlling delay jitter inside the network will also keep the rate jitter small; however, controlling rate jitter would not automatically eliminate delay jitter. From the network's point of view, the partial reconstruction of traffic patterns, or rate-jitter control inside the network, may be enough to offset the interactions between switches and the network load fluctuations; however, for some applications, a delay-jitter-free *service* is required [3, 4, 14]. If the network does not control delay jitter, applications may have to do buffering at the destination. For this reason, delay jitter control inside the network may be preferred; however, this consideration has to be balanced against the relative complexities of implementing delay-jitter control vs. rate-jitter control. One interesting observation is that both rate-jitter controlled and delay-jitter controlled service disciplines are non-work-conserving.

In order to offer performance guarantees, network resources need to be reserved at each switch for each connection. We consider three types of resources: buffer space, bandwidth and schedulability. The buffer space required for a connection to prevent packet loss can be calculated by multiplying the maximum packet arrival rate by the longest residence time of a packet in the switch. To provide deterministic guarantees [3], the overall bandwidth allocated to the connections on any link should not exceed the bandwidth of that link. Schedulability should also be considered. As pointed out in [6], scheduling saturation can occur even if bandwidth is not overbooked. For example, two packets with delay bounds of 4 time units and service time of 3 time units may arrive at a switch at the same time. In this situation, it is impossible to meet the delay bounds for both packets. To prevent network resources from being overbooked, a set of connection admission control conditions need to be tested before a new connection can be admitted. In general, buffer condition, bandwidth condition and schedulability condition need to be tested, as in the case of Jitter-EDD. Framing strategy used in Stop-and-Go and HRR ensures that scheduling saturation never occur if bandwidth is not overbooked, thus no additional schedulability test is needed in Stop-and-Go and HRR.

3 Rate-controlled Static-Priority Queueing

In this section, we will present a new service discipline called Rate-Controlled Static-Priority (RCSP) queueing. A RCSP server consists of two components: a rate-controller and a static priority scheduler. The rate-controller is responsible for correcting the distortion of traffic pattern introduced by the previous switch, and allocating bandwidth to different connections; the scheduler is responsible for allocating delays to different connections. The combination of a rate-controller and a static priority scheduler decouples the allocation of delay and bandwidth, and also simplifies admission control and implementation.

In this section, we discuss traffic characterization, describe the queueing algorithm and the corresponding admission control policies for RCSP, and finally, propose an implementation that is suitable for high-speed networks.

3.1 Traffic Specification

To provide performance guarantees, the network needs to allocate resources on a per-connection basis. This requires clients to specify their traffic characteristics. We use the (X_{min}, X_{ave}, I, P) traffic characterization as suggested in [6, 3]. Here X_{min} is the minimum packet inter-arrival time, X_{ave} is the average packet inter-arrival time over an averaging interval, I is the length of the interval, and P is the maximum packet size.

3.2 Algorithm

In this section, we will give a conceptual description of the RCSP queueing algorithm, the implementation issues are discussed in Section 3.4.

A rate-controlled static-priority server has two components: a rate controller and a static-priority scheduler. The rate controller shapes the input traffic from each connection into the desired traffic pattern by assigning an eligibility time to each packet; the scheduler orders the transmission of eligible packets from all the connections.

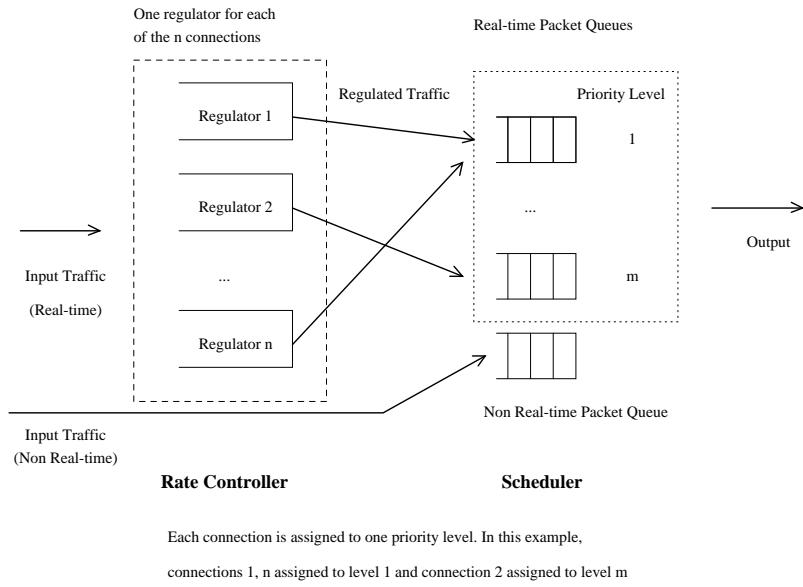


Figure 1: Rate-Controlled Static-Priority Queueing

Conceptually, a rate controller consists of a set of regulators corresponding to each of the connections traversing the switch; each regulator is responsible for shaping the input traffic of the corresponding connection into the desired traffic pattern. Regulators controls the interactions between switches and eliminates jitter. We present two types of regulators:

- *rate-jitter controlling regulator*, which controls rate jitter by partially reconstructing the traffic pattern;
- *delay-jitter controlling regulator*, which controls delay jitter by fully reconstructing the traffic pattern.

The regulator achieves this control by assigning each packet an eligibility time upon its arrival and holding the packet till that time before handing it to the scheduler. Different ways of calculating the eligibility time of a packet will result in different types of regulators.

For a rate-jitter controlling regulator, ET_j^k , the eligibility time of the k^{th} packet on a connection that traverses switch j , is defined with reference to the eligibility times of packets arriving earlier at the switch on the same connection:

- $ET_j^1 = AT_j^1$,
- $ET_j^k = \max(ET_j^{k-1} + X_{min} + \tau_j^k, AT_j^k)$, $k > 1$

where X_{min} is the minimum packet inter-arrival time for the connection, AT_j^k is the time the k^{th} packet on the connection arrived at switch j , and τ_j^k is the minimum number such that average packet inter-arrival time in any I period does not exceed X_{ave} . The computation of τ_j^k can be found in [13].

From this definition, we can see that $ET_j^k \geq AT_j^k$ always holds, i.e., a packet is never eligible before its arrival. Also, if we consider the sequence of packet eligibility times at switch j , $\{ET_j^k\}_{k=1,2,\dots}$, it always satisfies the (X_{min}, X_{ave}, I) traffic characterization.

For a delay-jitter controlling regulator, ET_j^k is defined with reference to the eligibility time of the same packet at the previous switch:

- $ET_0^k = AT_0^k$;
- $ET_j^k = ET_{j-1}^k + d_{j-1} + \pi_{j-1,j}$, $j > 1$,

where switch 0 is the source of the connection, d_{j-1} is the delay bound, or the maximum waiting time of packets on the same connection at the scheduler of switch $j-1$, and $\pi_{j-1,j}$ is the propagation delay between switch $j-1$ and switch j . (We will show in the next section how the waiting time in a scheduler can be bounded.)

From the second definition, we have:

- $ET_j^k \geq AT_j^k$, i.e., a packet is never eligible before its arrival;
- $ET_j^{k+1} - ET_j^k = AT_0^{k+1} - AT_0^k$, i.e., the traffic pattern of the connection at the *entrance* to the network is fully reconstructed at the output of the regulator of every switch; if the input traffic obeys the traffic specification at the entrance to the network, it will obey the traffic specification at the output of the regulator of each switch.

Both types of regulators will enforce the traffic specification requirement for each connection so that the traffic going into the scheduler will always satisfy the traffic specification.

The scheduler consists of a number of prioritized real-time packet queues and a non real-time packet queue (Figure 1). Packets at priority level 1 have the highest priority. A connection is assigned to a particular priority level at the connection's establishment time; all the packets from the connection will be inserted into the real-time packet queue at that priority level.

As will be discussed in section 3.3, there is a delay bound associated with each priority level; by limiting the number of connections at each priority level using the admission control conditions described in 3.3, the waiting time of each packet at a priority level is guaranteed to be less than the delay bound associated with that level.

The scheduler services packets using a non-preemptive static-priority discipline:

- when the server chooses the next packet to transmit, the packet at the head of the highest-priority non-empty real-time queue is always chosen;
- non real-time packets are transmitted only when there are no real-time packets in the scheduler; the order in which non real-time packets are served is not specified;
- the transmission of a lower-priority packet is not preempted by the arrival of a higher-priority packet.

- packets are always inserted at the end of the corresponding real-time queue, i.e., each real-time queue is served in FCFS order.

The algorithm described above would not provide performance guarantees without connection admission control conditions limiting the number of connections at each priority level. In the next section, we will describe the admission control conditions for RCSP. We will show that, if the admission control conditions are satisfied, by using the RCSP queueing described above, the end-to-end delay, delay jitter, throughput, and loss-free requirements of each admitted connection can be satisfied.

3.3 Admission Control Conditions

As described in section 1, two levels of control are needed to provide performance guarantees: service discipline at the packet level and connection admission control at the connection level. In this section, we will give the admission control conditions associated with the RCSP queueing discipline.

The residence time of a packet in a RCSP switch has two components: the *holding* time in the regulator and the *waiting* time in the scheduler (which includes both the queueing time and the service time). In the following, we first consider the case of a single switch; we show that, under certain admission control conditions, the waiting time of a packet is bounded by a fixed delay bound. We then examine the case of two switches connected in cascade; we show that if the connection admission control conditions are satisfied at both switches, the regulator of the second switch might delay a packet only if the packet is *ahead* of schedule in the first switch; the sum of the waiting time of a packet in the first switch and the holding time in the second switch is no greater than the delay bound of the first switch. From this result, we establish the end-to-end delay and delay jitter characteristics of a connection, and the buffer space requirement at each switch.

The following theorem gives the conditions to bound the waiting time of packets at each priority level in the scheduler.

Theorem 1 *Let d_1, d_2, \dots, d_n ($d_1 < d_2 < \dots < d_n$) be the delay bounds associated with each of the n priority levels, respectively, in an RCSP switch. Assume that the j^{th} connection among the i_k connections traversing the switch at priority level k has the traffic specification $(X_{\min k,j}, X_{\text{ave} k,j}, I_{k,j}, P_{k,j})$. Also assume that the link speed is l , and the size of the largest packet that can be transmitted onto the link is P_{\max} . If*

$$\sum_{k=1}^m \sum_{j=1}^{i_k} \lceil \frac{d_m}{X_{\min k,j}} \rceil P_{k,j} + P_{\max} \leq d_m l , \quad (1)$$

the waiting time of an eligible packet at level m is bounded by d_m .

Instead of giving a formal proof, we provide an intuitive basis for understanding the theorem. If inequality (1) holds, $\sum_{k=1}^m \sum_{j=1}^{i_k} \lceil \frac{d_m}{X_{\min k,j}} \rceil P_{k,j}$ is the maximum number of bits from priority levels 1 through m that can be transmitted by the scheduler without an idle interval or a lower level packet transmitted. The longest waiting time in the scheduler for a level- m packet corresponds to the case in which a lower-priority packet is being transmitted when the packet arrives at the scheduler, and is followed immediately by the longest possible transmission of packets with higher or equal priorities. Since the length of the lower-priority packet is bounded by P_{\max} , inequality (1) ensures that the level- m packet would be transmitted within d_m units of time after its arrival even in the worst case.

Inequality (1) in Theorem 1 gives the admission control condition for switch k . For each priority level m , a state variable $\text{LongestBusyPeriod}_m$ is kept and is initialized to be P_{\max} . When an establishment request for a new connection with traffic specification $(X_{\min}, X_{\text{ave}}, I, P)$ comes in, the following tests are performed from priority level 1 to n :

$$\text{LongestBusyPeriod}_m + \lceil \frac{d_m}{X_{\min}} \rceil P \leq d_m l \quad (2)$$

The new connection can be placed into priority level k without jeopardizing the performance guarantees of other existing connections if only if equality (2) holds for $k, k+1, \dots, n$. Usually, there are multiple priority levels that the new connection can be placed into. Different connection establishment schemes can be used to choose the priority level for the new connection. One scheme of establishing real-time connection is described in [6], in which a new connection is setup with a one-round trip message: during the forward trip of the establishment message, the new connection is assigned the smallest possible delay bound at each switch to maximize the probability that the new connection's end-to-end delay bound can be met; during the return trip, the local delay bound at each switch can be relaxed if it turns out that the offered delay bound is tighter than the requested one. This is a greedy delay allocation policy: although it minimizes the blocking probability of the new connection, it increases the blocking probability for other connection establishment requests before the delay bound for the first connection is relaxed. How to assign the local delay bound, or priority level, to a new connection request during the forward trip is a *policy* issue, which deserves further research. If we decide to place the new connection in priority level h , state variables need to be updated accordingly:

$$\text{LongestBusyPeriod}_m \leftarrow \text{LongestBusyPeriod}_m + \lceil \frac{d_m}{X_{\min}} \rceil P \quad m = h, \dots, n$$

The admission control conditions in theorem 1 avoid both bandwidth and scheduling saturation. The tests have just $O(n)$ steps, where n is the number of priority levels, and are simple enough to be used in some fast connection establishment scheme [12].

Theorem 1 just gives the bounds for the waiting times of packets in the scheduler; the holding times of packets in the regulators have not been taken into account. In the previous section, we have introduced two types of regulators: the rate-jitter controlling regulator and the delay-jitter controlling regulator. In the following lemma, we establish the delay characteristics of two switches connected in cascade. when each type of regulator is used.

Lemma 1 *Consider a connection with traffic specification $(X_{\min}, X_{\text{ave}}, I, P)$ passing through two RCSP switches connected in cascade. Assume that the connection is assigned to the priority level with delay bound d at the first switch. Furthermore, assume that the admission control policy as defined in Theorem 1 is satisfied at both switches. It follows that, for any packet on the connection,*

- the waiting time in the first switch plus the holding time in the second switch $\leq d$ if a rate-jitter controlling regulator is used in the second switch;
- the waiting time at the first switch plus the holding time at the second switch $= d$ if a delay-jitter controlling regulator is used at the second switch.

The key observation for understanding the lemma is that a packet is held in a regulator only when the packet was transmitted ahead of schedule by the previous switch; the amount of holding time in the regulator is never greater than the amount of time the packet is ahead of schedule. When a rate-jitter controlling regulator is used, the amount of time a packet is to be held is calculated according to the packet spacing requirement, which may be less than the amount of time it was ahead of schedule at the previous switch. When a delay-jitter controlling regulator is used, the amount of holding time is exactly the amount of time the packet was ahead of schedule when it left the previous scheduler.

The following theorem establish the conditions for bounding end-to-end delay for a connection.

Theorem 2 *Consider a connection with traffic specification $(X_{\min}, X_{\text{ave}}, I, P)$ passing through k switches connected in cascade where the end-to-end propagation delay is Π . Assume that the connection is assigned to the priority levels with delay bounds d_{i_1}, \dots, d_{i_k} at each switch, respectively. Furthermore, assume that the admission control policy defined in Theorem 1 is satisfied at all the switches. If the traffic on the connection obeys the $(X_{\min}, X_{\text{ave}}, I, P)$ specification at the entrance to the first switch,*

1. the end-to-end delay for any packet on the connection is bounded by $\sum_{j=1}^k d_{i_j} + \Pi$ if rate-jitter controlling regulators are used;
2. the end-to-end delay and delay jitter for any packet are bounded by $\sum_{j=1}^k d_{i_j} + \Pi$ and d_{i_k} respectively if delay-jitter controlling regulators are used;
3. an amount of buffer space equal to $(\lceil \frac{d_{i_{j-1}}}{X_{min}} \rceil + \lceil \frac{d_{i_j}}{X_{min}} \rceil)P$ is needed at switch j to prevent packet loss
 $(j = 1, \dots, k; d_{i_0} = 0)$.

For the first two parts of the theorem, consider the end-to-end delay of one packet in the connection, which can be expressed as:

$$\sum_{j=1}^k (h_j + w_j) + \Pi$$

where h_j and w_j are the holding and the waiting times at switch j . If we re-arrange the terms, the formula becomes:

$$h_1 + \sum_{j=1}^{k-1} (w_j + h_{j+1}) + w_k + \Pi$$

If the traffic obeys the (X_{min}, X_{ave}, I, P) characterization at the entrance to the first switch, there is no holding time in the first regulator, so $h_1 = 0$. Also from Lemma 1, $w_j + h_{j+1} = d_{i_j}$ holds for a delay-jitter controlling regulator and $w_j + h_{j+1} \leq d_{i_j}$ holds for a rate-jitter controlling regulator; furthermore, $w_k \leq d_{i_k}$ holds in both cases; the correctness of the first two parts of the theorem immediately follows.

To understand the third part of the theorem, notice that the longest times a packet can stay in the regulator and the scheduler of the j_{th} switch are $d_{i_{j-1}}$ and d_{i_j} , respectively; since the minimum packet inter-arrival time is X_{min} , it follows that the maximum numbers of packets at the regulator and the scheduler are $\lceil \frac{d_{i_{j-1}}}{X_{min}} \rceil$ and $\lceil \frac{d_{i_j}}{X_{min}} \rceil$, respectively.

3.4 Implementation

In this section, we present one implementation of the RCSP queueing discipline. We believe that this implementation is simple enough to run at very high speed.

We have shown in Section 3.2 that a RCSP server has two components, a scheduler and a rate controller. The scheduler consists of multiple prioritized FCFS queues, and the rate controller consists of a set of regulators corresponding to each connection. We implement each real-time packet queue in the scheduler as a linked list, and the rate controller as a modified version of a calendar queue [1, 13]. Notice that the conceptual decomposition of the rate controller into a set of regulators does not imply that there are multiple physical regulators in an implementation; a common mechanism can be shared by all regulators. Each regulator has two functions: computing the eligibility times for incoming packets on the corresponding connection, and holding packets till their eligible times. Eligibility times for packets from different connections are computed using the same formula (either the formula for the rate-jitter controlling regulator, or the formula for the delay-jitter controlling regulator as described in Section 3.2) with different parameters; holding packets is equivalent to managing a set of timers; the mechanism for managing timers, which is a calendar queue, can be shared by all regulators.

Figure 2 shows the proposed implementation. Each of the real-time queues is implemented as a linked list. The operations performed on a linked list are: deleting the packet at the head (when the next packet is chosen for transmission) and appending a packet or linked list of packets at the tail (when a packet arrives or the clock ticks; see below).

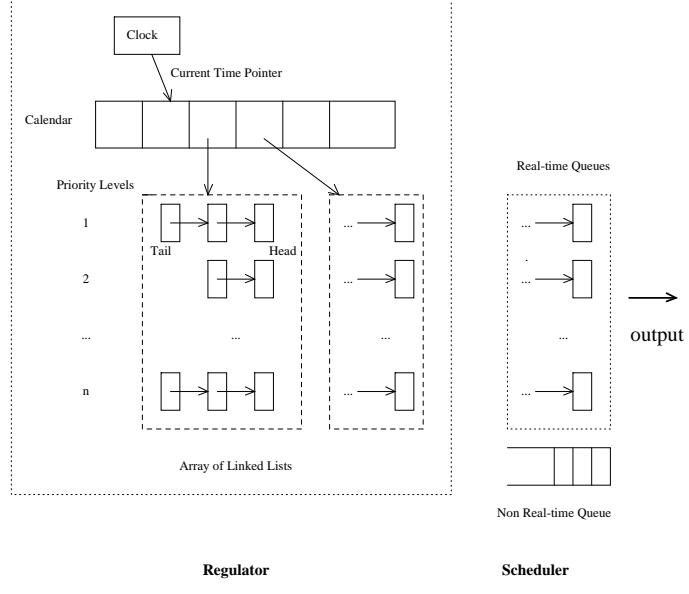


Figure 2: Rate-controlled Multi-priority Queueing

The rate controller is implemented using a modified version of a calendar queue. A calendar queue consists of a clock and a calendar, which is a pointer array indexed by time; each entry in the calendar points to an array of linked lists indexed by priority levels.

The clock ticks at fixed time intervals. Upon every tick of the clock, the linked lists in the array indexed by the current time are appended at the end of the scheduler's linked lists: packets from the linked list of one priority level in the rate-controller are appended to the linked list of the same priority level in the scheduler.

Upon the arrival of each packet, the eligibility time of the packet, ET , is calculated; if $\lfloor \frac{ET}{Tick} \rfloor$ is equal to the current clock time, where $Tick$ is the clock tick interval, the packet is appended at the end of the corresponding real-time queue of the scheduler; otherwise, the packet is appended at the corresponding linked list at the calendar queue entry indexed by $\lfloor \frac{ET}{Tick} \rfloor$.

Notice that in this implementation, packets may become eligible upto $Tick$ time units earlier than their actual eligibility times. The admission control conditions need to be modified to accommodate for this; the inequality (1) in theorem 1 should now be:

$$\sum_{k=1}^m \sum_{j=1}^{i_k} \lceil \frac{d_k + Tick}{X_{\min k,j}} \rceil P_{k,j} + P_{\max} \leq d_{ml}$$

$d_k + Tick$ is used instead of d_k in the numerator to account for rounding errors in eligibility time due to the use of a calendar queue. Similar changes need to be made to Theorem 2: the buffer space requirement is now $(\lceil \frac{d_{i,j-1} + Tick}{X_{\min}} \rceil + \lceil \frac{d_{i,j}}{X_{\min}} \rceil)P$; for a delay-jitter controlled connection, the end-to-end delay jitter bound is now $d_{i_k} + Tick$. A smaller $Tick$ would introduce less rounding error, but would have a larger overhead. The effect of $Tick$ size on performance should be studied further.

As can be seen, the data structures used in the proposed implementation are simple: arrays and linked lists; the operations are all constant-time ones: insertion at the tail of a linked list and deletion from the head of a linked list. We believe that it is feasible to implement this in a very high speed switch.

4 Comparison and Analysis

In this section, we compare the RCSP service discipline with the previously proposed non-work-conserving disciplines: Jitter-EDD, Stop-and-Go, and HRR. The goal of this comparison is to determine which of the desired features are provided by each discipline and to identify the linkage between *mechanisms* and *features*. We investigate the common mechanisms in these disciplines that contribute to such desired features as bounded delays, bounded delay jitters, decoupling of delay and bandwidth allocation, evenly distributed buffer space reservation to prevent packet loss, and easy implementation in high-speed switches.

Desired feature	Mechanism	HRR	SG	J-EDD	RCSP
bounded delay	sorted priority queue			✓	
	static multi-level priorities	✓	✓		✓
bounded delay jitter	reconstruction of traffic patterns at each switch		✓	✓	✓
decoupled delay and bandwidth allocation	use two components: regulator and scheduler			✓	✓
uniformly distributed buffer space	shaping traffic pattern into desired form at each switch	✓	✓	✓	✓

Table 2: Mechanisms vs. Desired Features in Service Disciplines for Real-time Communication

In table 2, we can see that each desire feature is implemented by nearly identical mechanism in each service discipline. The only exception is that Jitter-EDD uses sorted priority queue to allocate delays, while HRR, Stop-and-Go and RCSP use static multi-level priorities. Sorted priority queue allows the maximum flexibility in allocating delays: there is a continuous spectrum of priorities that can be offered to connections; however, it may be difficult to implement sorted priority queues at high speed; also, a complicated schedulability test is required at connection establishment time. A static multi-level priority scheme is not as flexible as a sorted priority queue, in the sense that it can only allocate a fixed number of priorities to connections; however, it has the advantages of simplicity in admission control tests and in implementation; also, in practice, it is usually sufficient to offer a fixed number of priorities to different connections.

HRR and Stop-and-Go use multi-level framing to perform both delay control and rate control; this introduces dependencies between delay and bandwidth allocation. The separation of rate control and delay control in Jitter-EDD and RCSP eliminates such dependencies. The regulator does not only allocate bandwidth, but also isolates the interactions between adjacent switches, and controls jitter. The scheduler is responsible for servicing packets with different priorities to meet a variety of delay requirements. Following this separation and having different combinations of regulators and schedulers will result in a class of service disciplines as shown in Figure 3.

The combination of rate-jitter controlling regulators or delay-jitter controlling regulators with a static priority scheduler is the proposed RCSP discipline; the combination of delay-jitter controlling regulators and an Earliest-Due-Date scheduler is the Jitter-EDD discipline. A service discipline with rate-jitter controlling regulators or delay-jitter controlling regulators and a FCFS scheduler has been proposed in [13] for its simplicity; however, it may not be able to provide different delay bounds to different connections.

Controlling traffic pattern distortions inside the network is used in four disciplines to achieve efficient buffer allocation. From the network's point of view, rate-jitter control or partial reconstruction of traffic patterns inside the network is sufficient to control the interaction between switches and the network load fluctuations. However, delay-jitter control may be preferred in the cases when a delay jitter bounded service is desired [3, 4, 14]. The important consideration here regards the relative

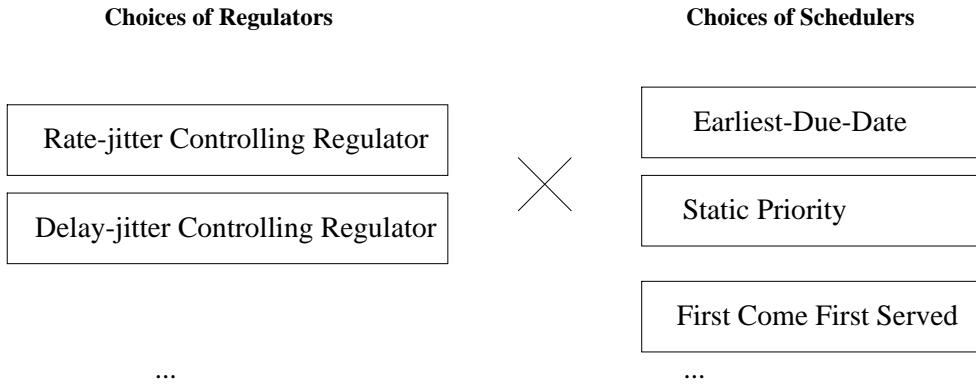


Figure 3: Servers consist of Regulator and Scheduler

complexities of implementing delay-jitter control vs. rate-jitter control. In the two framing-based solutions, Stop-and-Go needs to synchronize between input links and output links at each switch in order to achieve delay-jitter control; HRR does not have this synchronization, and thus can not bound delay-jitter in its current form; however, it is possible to modify HRR to implement Stop-and-Go queueing by framing and synchronizing input and output links [10]. For a service discipline with a rate controller, like Jitter-EDD and RCSP, a calendar queue mechanism can be used to implement both the rate-jitter controlling and the delay-jitter controlling regulators. In this case, the only difference between a delay-jitter controlling regulator and a rate-jitter controlling regulator is the computation of the eligibility time. A delay-jitter controlling regulator needs one more quantity to compute the eligibility time of a packet than a rate-jitter controlling regulator: “the amount of time the packet was ahead of schedule in the previous switch”. This quantity can be calculated in the previous switch and stamped into the packet header as proposed in [14]. However, this incurs the overhead of timestamping and transmission of the timestamp. This overhead may be very significant in an ATM network, where the cell size is just 53 bytes; however, it may not be a big problem in a fast packet networks where packet size is in the order of kilobytes. So, for an ATM network, a rate-jitter controlling mechanism may be more cost effective; for a fast packet network, a delay-jitter controlling mechanism can provide an almost “free” delay-jitter-bounded service.

5 Conclusions and Future Work

We have proposed a new service discipline, called the Rate-Controlled Static-Priority (RCSP) queueing discipline, that can provide throughput, delay, delay jitter, and loss free guarantees in a connection-oriented packet-switching network. We have given the corresponding admission control conditions, and proposed an implementation that is suitable for high-speed networks. RCSP achieves flexibility in the allocation of delay and bandwidth as well as simplicity of admission control and implementation. The key idea is to separate rate-control and delay-control functions in the design of the server. By having different combinations of rate controllers and schedulers, we get a class of service disciplines, of which RCSP and the previously proposed Jitter-EDD are two instances.

There are still several important issues left to be investigated, some of which are the following:

- Unlike Jitter-EDD, which can offer a continuous spectrum of delay bounds, RCSP can only offer a fixed number of delay bounds. The tradeoffs in choosing the number of priority levels and the delay bound associated with each priority level are not fully understood and need to be studied further.
- The proposed implementation uses a calendar queue mechanism to implement the rate-controller.

A smaller clock tick in the calendar queue would result in a better performance, but would also have a larger overhead. A simulation study is under way to study the effects of the clock tick size on the performance.

- Another issue is the relationship between non-work-conserving disciplines and the statistical multiplexing principles of packet switching networks. The fundamental difference between packet switching and circuit switching is that packet switching allows statistical multiplexing. Statistical multiplexing introduces both flexibility of bandwidth sharing and interferences between different conversations. Non-work-conserving disciplines limit the effect on different conversations due to network fluctuations by *holding* packets. One potential argument against this holding is that it may reduce the flexibility offered by statistical multiplexing. This problem should be considered in the context of a larger question, which is unfortunately an open one, how to achieve both the flexibility and guaranteed performance in packet-switching networks.
- A related issue is how to enhance the utilization of the network by offering statistical guarantees in RCSP. Solutions have been proposed for Delay-EDD [6] and Stop-and-Go [8] to overbook resources by exploiting the traffic burstiness. We are currently extending those ideas to RCSP.

References

- [1] Randy Brown. Calendar queues: A fast O(1) priority queue implementation for the simulation event set problem. *Communications of the ACM*, 31(10):1220–1227, October 1988.
- [2] Alan Demers, Srinivasan Keshav, and Scott Shenker. Analysis and simulation of a fair queueing algorithm. In *Journal of Internetworking Research and Experience*, pages 3–26, October 1990. Also in Proc. ACM SIGCOMM’89, pp 3-12.
- [3] Domenico Ferrari. Client requirements for real-time communication services. *IEEE Communications Magazine*, 28(11), November 1990. also RFC 1193.
- [4] Domenico Ferrari. Design and applications of a delay jitter control scheme for packet-switching internetworks. In *Proceedings of the Second International Workshop on Network and Operating System Support for Digital Audio and Video*, Heidelberg, Germany, 1991.
- [5] Domenico Ferrari. Real-time communication in an internetwork. Technical Report TR-91-001, International Computer Science Institute, Berkeley, California, January 1992.
- [6] Domenico Ferrari and Dinesh Verma. A scheme for real-time channel establishment in wide-area networks. *IEEE Journal on Selected Areas in Communications*, 8(3):368–379, April 1990.
- [7] S. Jamaloddin Golestani. A stop-and-go queueing framework for congestion management. In *SIGCOMM’90 Symposium, Communications Architecture & Protocols*, pages 8–18, Philadelphia Pennsylvania, September 1990. ACM SIGCOMM.
- [8] S. Jamaloddin Golestani. Duration-limited statistical multiplexing of delay-sensitive traffic in packet networks, April 1991.
- [9] Charles R. Kalmanek, Hemant Kanakia, and Srinivasan Keshav. Rate controlled servers for very high-speed networks. In *IEEE Global Telecommunications Conference*, San Diego, California, December 1990.
- [10] Srinivasan Keshav, October 1991. personal communication.

- [11] Donald E. Knuth. *The Art of Computer Programming. Volume 3: Sorting and searching.* Addison-Wesley, 1975.
- [12] P. Todorova and Dinesh Verma. Fast establishment of real-time channels. Technical Report TR-89-056, International Computer Science Institute, Berkeley, California, October 1989.
- [13] Dinesh Verma. *Guaranteed Performance Communication in High Speed Networks.* PhD dissertation, University of California at Berkeley, November 1991.
- [14] Dinesh Verma, Hui Zhang, and Domenico Ferrari. Guaranteeing delay jitter bounds in packet switching networks. In *Proceedings of Tricomm'91*, Chapel Hill, North Carolina, April 1991.
- [15] Hui Zhang and Keshav Srinivasan. Comparison of rate-based service disciplines. In *SIGCOMM'91 Conference, Communications Architecture & Protocols*, pages 113–122, Zurich, Switzerland, September 1991. ACM SIGCOMM.
- [16] Lixia Zhang. Virtual clock: A new traffic control algorithm for packet switching networks. In *SIGCOMM'90 Symposium, Communications Architecture & Protocols*, pages 19–29, Philadelphia Pennsylvania, September 1990. ACM SIGCOMM.