

Determining end-to-end delay bounds in heterogeneous networks

Pawan Goyal, Simon S. Lam, Harrick M. Vin

Department of Computer Sciences, University of Texas at Austin, Taylor Hall 2.124, Austin, Texas 78712-1188, USA
e-mail: {pawang,lam,vin}@cs.utexas.edu

Abstract. We define a class of Guaranteed Rate (GR) scheduling algorithms. The GR class includes Virtual Clock, Packet-by-Packet Generalized Processor Sharing and Self-Clocked Fair Queuing. For networks that employ scheduling algorithms belonging to GR, we present a method for determining an upper bound on end-to-end delay. The method facilitates determination of end-to-end delay bounds for a variety of sources. We illustrate the method by determining end-to-end delay bounds for sources conforming to Leaky Bucket and exponentially bounded burstiness.

1 Introduction

Computer networks have advanced to a point where they can support multimedia applications like audio and video conferencing and multimedia information retrieval. Such applications require the network to provide a wide range of quality-of-service (QoS) guarantees (including minimum bandwidth, packet delay, delay (jitter and loss)). Whereas the minimum guaranteed bandwidth must be large enough to accommodate motion video of acceptable resolution, the end-to-end delay must be small enough for interactive communication. In order to avoid breaks in continuity of audio and video playback, delay jitter and loss must be sufficiently small. Techniques for determining an upper bound on end-to-end delay in a network is the subject matter of this paper.

The end-to-end delay of a packet depends on the source traffic characteristics and the scheduling algorithm at the network switches. In the recent past, several source specifications including Leaky Bucket [13], exponentially bounded burstiness (EBB) [16], Flow Specification [11] and the Tenet model [6] have been studied. The scheduling algorithms that have been proposed include Stop and Go Queuing [8], Delay EDD [18], Jitter EDD [18], Hierarchical Round Robin [10], Rate Control Static Priority Queuing [17], Self-Clocked Fair Queuing (SCFQ) [9], Virtual Clock [19] and Packet-by-Packet Generalized Processor Sharing (PGPS) [13]. A survey of the scheduling algorithms can be found in [18]. The problem of determining an upper bound on end-to-end delay has also received considerable attention [2, 3, 11, 13,

16]. However, most of the techniques determine end-to-end delay by considering a specific source traffic specification and scheduling algorithm. In an integrated network supporting audio, video and data services, the sources have widely varying characteristics. Moreover, in a wide area networking environment, each switch may employ a different scheduling algorithm. Methods for determining end-to-end delay of packets in such heterogeneous environments have not received much attention.

In this paper, we take a step towards addressing the above limitation by (1) defining a class of *Guaranteed Rate (GR)* scheduling algorithms, and (2) developing a method for determining an upper bound on end-to-end delays for a network of switches, each of which employs a scheduling algorithm in the GR class. We demonstrate that many of the scheduling algorithms proposed in the literature (e.g. Virtual Clock, PGPS and SCFQ) belong to the class of GR scheduling algorithms. We also show that the method for determining an upper bound on end-to-end delay is general and can be used to determine delay bounds for various source traffic specifications. We employ the method to derive a deterministic end-to-end delay bound for Leaky Bucket sources (a deterministic source characterization) and an upper bound on the tail distribution of the delay for EBB sources (a stochastic characterization). The end-to-end delay bounds that we derive are parametrized by the scheduling algorithm used at each switch and can be instantiated to derive delay bounds for a specific scheduling algorithm or a combination of scheduling algorithms.

The rest of the paper is organized as follows. In Sect. 2, we define the class of GR scheduling algorithms. In Sect. 3, we present a method for determining end-to-end delay bounds and finally, Sect. 4 summarizes our results.

2 Guaranteed Rate scheduling algorithms

Each unit of data transmission at the network level is a packet. We refer to the sequence of packets transmitted by a source as a *flow* [19]. Each packet within a flow is serviced by a sequence of servers (switching elements) along the path from the source to the destination in the network. To provide guaranteed performance, the servers reserve a rate for a flow

and employ a rate-based scheduling algorithm. Based on this rate reservation, many scheduling algorithms can guarantee a deadline by which a packet of a flow will be transmitted. This guarantee is referred to as *delay guarantee*. We refer to the class of rate-based scheduling algorithms which provide such guarantees as *Guaranteed Rate (GR)* scheduling algorithms.

The delay guarantees provided by these algorithms are based on the Guaranteed Rate *clock*¹ values associated with each packet. To define the Guaranteed Rate clock values, consider a flow f that is associated with rate r_f (in bits/s). Let p_f^j and l_f^j denote the j^{th} packet of flow f and its length, respectively. Additionally, let $GRC^i(p_f^j)$ and $A^i(p_f^j)$ denote the Guaranteed Rate clock value and arrival time of packet p_f^j , at server i , respectively. Then, GR clock value for a packet is given by:

$$GRC^i(p_f^0) = 0 \quad (1)$$

$$GRC^i(p_f^j) = \max\{A^i(p_f^j), GRC^i(p_f^{j-1})\} + \frac{l_f^j}{r_f} \quad (2)$$

We use the Guaranteed Rate clock value of a packet to define the class of GR scheduling algorithms as follows.

Definition 1. A scheduling algorithm at server i belongs to class GR for flow f if it guarantees that packet will be transmitted by $GRC^i(p_f^j) + \beta^i$, where β^i is a constant which depends on the scheduling algorithm and the server.

The concept of delay guarantee to a packet based on its expected arrival time was introduced in [11, 15]. As is evident from the definition (also observed in [11]), a key property of the class of GR scheduling algorithms is that they provide a delay guarantee for a source independent of the behavior of other sources in the network, and hence isolate the sources. Isolation of sources has been considered to be a very desirable property of scheduling algorithms, especially in large heterogeneous networks where sources may be malicious [1, 5, 11, 14]. This property is the basis for development of a conceptually simple method for determining end-to-end delay bounds in Sect. 3.

In the next few subsections we demonstrate that the GR class includes Virtual Clock, PGPS and SCFQ with certain rate assignments.

2.1 Virtual clock

The Virtual Clock scheduling algorithm assigns each packet a virtual clock value on its arrival and orders the transmission of packets by increasing virtual clock values. If flow f is assigned rate r_f at server i , then the virtual clock value for packet of flow f at server i , denoted by $VC^i(p_f^j)$, is computed as follows [19]:

$$VC^i(p_f^0) = 0 \quad (3)$$

¹ The concept of guaranteed rate clock is the same as the concept of virtual clock in [19]. We coin a new term for the virtual clock concept to avoid any confusion between the general concept of virtual clock and as it relates to the virtual clock scheduling algorithm.

$$VC^i(p_f^j) = \max\{A^i(p_f^j), VC^i(p_f^{j-1})\} + \frac{l_f^j}{r_f} \quad (4)$$

The following delay guarantee was presented in [15]. Define a flow to be active at time t if $VC^i(p_f^j) \geq t$, where p_f^j is the last packet of flow f that has arrived before time t . Let $a^i(t)$ denote the set of active flows f at server i at time t . Server i with capacity C^i is defined to have exceeded its capacity at time t if $\sum_{n \in a^i(t)} r_n > C^i$. Let l_{\max}^i be the maximum length of the packet served by server i and $L_{VC}^i(p_f^j)$ be the time at which the transmission of packet p_f^j is completed. Then, if a server's capacity has not been exceeded and flow f is assigned rate r_f ,

$$L_{VC}^i(p_f^j) \leq VC^i(p_f^j) + \frac{l_{\max}^i}{C^i} \quad (5)$$

Since the equations for virtual clock and GR clock are the same, it can be easily observed that if rate r_f is assigned to flow then the Virtual Clock scheduling algorithm belongs to GR for flow f with $\beta^i = \frac{l_{\max}^i}{C^i}$.

2.2 Packet-by-packet generalized processor sharing

The PGPS scheduling algorithm is a practical realization of generalized processor sharing (GPS) service discipline [13]. In fact, PGPS simulates GPS such that

$$L_{PGPS}^i(p_f^j) \leq L_{GPS}^i(p_f^j) + \frac{l_{\max}^i}{C^i} \quad (6)$$

where $L_{PGPS}^i(p_f^j)$ and $L_{GPS}^i(p_f^j)$ are the times at which packet p_f^j leaves server i employing the PGPS scheduling algorithm and the GPS service discipline, respectively, and l_{\max}^i is the maximum length of the packet served by server i [13]. Hence, to show that PGPS belongs to the class of GR scheduling algorithms, we establish a relationship between $L_{GPS}^i(p_f^j)$ and $GRC^i(p_f^j)$.

In GPS, each flow f is associated with a constant ϕ_f^i at server i . GPS is defined such that if flow f is backlogged at time t , it receives service at the rate of $\frac{\phi_f^i C^i}{\sum_{k \in b^i(t)} \phi_k^i}$, where C^i is the capacity of the server and $b^i(t)$ is the set of backlogged flows at GPS server i at time t . Thus, a GPS server is defined to have assigned rate r_f to flow f if $\frac{\phi_f^i}{\sum_{k \in b^i(t)} \phi_k^i} \geq r_f$ whenever flow f is backlogged. Hence, if flow f has been assigned rate r_f then

$$L_{GPS}^i(p_f^j) \leq \max\{A^i(p_f^j), L_{GPS}^i(p_f^{j-1})\} + \frac{l_f^j}{r_f} \quad j \geq 1 \quad (7)$$

Let $L_{GPS}^i GRC^i(p_f^0) = GRC^i(p_f^0) = 0$. From (7) and (2), it can be easily shown that

$$L_{GPS}^i(p_f^j) \leq GRC^i(p_f^j) \quad j \geq 1 \quad (8)$$

From (8) and (6), we get

$$L_{PGPS}^i(p_f^j) \leq GRC^i(p_f^j) + \frac{l_m^{\max}}{C^i} \quad (9)$$

Hence, if rate r_f is assigned to flow f at server i , then the PGPS scheduling algorithm belongs to GR for flow f with $\beta^i = \frac{l_m^{\max}}{C^i}$.

2.3 Self-Clocked Fair Queuing

The SCFQ scheme, originally proposed in [4] and later analyzed in [9], was designed to facilitate the implementation of a fair queuing scheme. The scheme is defined as follows.

1. On arrival, a packet p_f^j is stamped with service tag $F^i(p_f^j)$, computed as:

$$F^i(p_f^0) = 0 \quad (10)$$

$$F^i(p_f^j) = \max\{F^i(p_f^{j-1}), v^i(A^i(p_f^j))\} + \frac{l_f^j}{\phi_f^i} \quad j \geq 1 \quad (11)$$

where ϕ_f^i is a constant associated with flow f at server i .

2. The server virtual time at time t , $v^i(t)$, is defined to be equal to the service tag of the packet being service at time t . $v^i(t) = t$ when the server is idle.

3. Packets are serviced in increasing order of their service tags.

The following theorem proves that SCFQ algorithm also belongs to the class of GR scheduling algorithms. Let c^i be the set of flows serviced by server i employing SCFQ scheduling algorithm.

Theorem 1. *If $\sum_{m \in c^i} r_m \leq C^i$ and $\forall m \in c^i : \phi_m^i = r_m$, then the departure time of packet p_f^j in SCFQ (denoted by $L_{SCFQ}^i(p_f^j)$) is given by*

$$L_{SCFQ}^i(p_f^j) \leq GRC^i(p_f^j) + \sum_{m \in c^i \wedge m \neq f} \frac{l_m^{\max}}{C^i} \quad (12)$$

where l_m^{\max} is the maximum length for packets in flow m .

Proof. Let set B_f^i be defined as follows.

$$B_f^i = \{n \mid n > 0 \wedge F^i(p_f^{n-1}) < v^i(A^i(p_f^n))\}$$

Let $k \leq j$ be the largest integer in B_f^i . Since packets of a flow are serviced in FCFS order, they form a queue. Let R_m be the set of packets of flow m , $m \neq f$, served in the interval $[A^i(p_f^k), L_{SCFQ}^i(p_f^j)]$, excluding the packet that was at the head of its flow queue at time $A^i(p_f^k)$. Let $v_m(t)$ be the finish time of the packet at the head of the queue of flow m at time t . If no packet is present, then set $v_m(t) = v^i(t)$. Since SCFQ services packets in increasing order of the service tags,

$$v_m(A^i(p_f^k)) + \sum_{n \in R_m} \frac{l_m^n}{r_m} \leq v^i(A^i(p_f^k)) + \sum_{n=0}^{n=j-k} \frac{l_f^{k+n}}{r_f} \quad (13)$$

Since $v_m(A^i(p_f^k)) \geq v^i(A^i(p_f^k))$ we have

$$\sum_{n \in R_m} \frac{l_m^n}{r_m} \leq \sum_{n=0}^{n=j-k} \frac{l_f^{k+n}}{r_f} \quad (14)$$

$$\Rightarrow \sum_{n \in R_m} l_m^n \leq \sum_{n=0}^{n=j-k} \frac{l_f^{k+n}}{r_f} r_m \quad (15)$$

By adding (15) over all flows other than f , we have

$$\sum_{m \in c^i \wedge m \neq f} \sum_{n \in R_m} l_m^n \leq \sum_{n=0}^{n=j-k} \frac{l_f^{k+n}}{r_f} \sum_{m \in c^i \wedge m \neq f} r_m \quad (16)$$

Adding $\sum_{n=0}^{n=j-k} l_f^{k+n}$ to both sides of the equation

$$\sum_{m \in c^i \wedge m \neq f} \sum_{n \in R_m} l_m^n + \sum_{n=0}^{n=j-k} l_f^{k+n} \leq \sum_{n=0}^{n=j-k} \frac{l_f^{k+n}}{r_f} \sum_{m \in c^i} r_m \quad (17)$$

$$\begin{aligned} &\Rightarrow \frac{\sum_{m \in c^i \wedge m \neq f} \sum_{n \in R_m} l_m^n + \sum_{n=0}^{n=j-k} l_f^{k+n}}{C^i} \\ &\leq \frac{\sum_{n=0}^{n=j-k} \frac{l_f^{k+n}}{r_f} \sum_{m \in c^i} r_m}{C^i} \end{aligned} \quad (18)$$

Since $\sum_{m \in c^i} r_m \leq C^i$,

$$\begin{aligned} &\frac{\sum_{m \in c^i \wedge m \neq f} \sum_{n \in R_m} l_m^n + \sum_{n=0}^{n=j-k} l_f^{k+n}}{C^i} \\ &\leq \sum_{n=0}^{n=j-k} \frac{l_f^{k+n}}{r_f} \end{aligned} \quad (19)$$

The left-hand side of (19), denoted by T is the time spent in servicing the packets of flow f and packets of flows $m (m \neq f)$, excluding the packet which was at the head of the flow queue at time $A^i(p_f^k)$, before p_f^j departs. Since the server may spend some time servicing the packets that were at the head of the flow queues, we have

$$A^i(p_f^k) + \sum_{m \in c^i \wedge m \neq f} \frac{l_m^{\max}}{C^i} + T \geq L_{SCFQ}^i(p_f^j) \quad (20)$$

where l_m^{\max} is the maximum length of flow m packet. From (19), we get

$$\begin{aligned} &A^i(p_f^k) + \sum_{m \in c^i \wedge m \neq f} \frac{l_m^{\max}}{C^i} + \sum_{n=0}^{n=j-k} \frac{l_f^{k+n}}{r_f} \\ &\geq L_{SCFQ}^i(p_f^j) \end{aligned} \quad (21)$$

From (2) we get

$$GRC^i(p_f^j) + \sum_{m \in c^i \wedge m \neq f} \frac{l_m^{\max}}{C^i} \geq L_{SCFQ}^i(p_f^j) \quad (22)$$

Hence, if $\sum_{m \in c^i} r_m \leq C^i$ and $\forall m \in c^i : \phi_m^i = r_m$, then the SCFQ scheduling algorithm belongs to GR for flow f with $\beta^i = \sum_{m \in c^i \wedge m \neq f} \frac{l_m^{\max}}{C^i}$. Note also that Theorem 1 provides a tight bound over all rate assignments. An arrival sequence in which the upper bound is realized can be constructed from the proof steps.

3 Determining end-to-end delay bounds

3.1 Method

A simple method for determining a bound on end-to-end packet delays is to consider each server in isolation, and compute the summation of the maximum delay at each server along the path from the source to the destination [2, 3]. Though this may be the only feasible approach for many scheduling algorithms, it has the following drawbacks:

- Due to the variability in the delay experienced by packets at a server, the shape of the traffic becomes distorted as it traverses through the network. Therefore, even if a source traffic specification is known at the first server on the path of the flow, it is difficult to determine the specification at a server further down on the path. This makes determining the delay at each of the server difficult.
- In many service disciplines, if a packet experiences a high delay at a server, it may experience a lower delay at the next server along the path. If each server is considered in isolation, the dependence between the delay experienced by packets at different servers is not accounted for, and hence the bound on the delay may be very conservative.

These limitations have been addressed by analyzing the path as a whole for a network of servers employing PGPS and sources conforming to Leaky Bucket [13]. Similarly, the approach has been utilized for burst scheduling networks with sources conforming to Flow Specification [11]. In this section, we generalize the approach to a heterogeneous network of servers each of which employs a scheduling algorithm in GR for any source specification. We derive a delay guarantee for a network of servers and reduce the problem of determining end-to-end delay to that of determining delay at a single server.

Let K be the total number of servers along the path of a flow, and let the i^{th} server on the path be denoted by i . Also, let server 0 be the source and server $K + 1$ be the destination. Since server K guarantees that packet (p_f^j) will be transmitted by $GRC^K(p_f^j) + \beta^K$ and the packet arrives at the first server at time $A^1(p_f^j)$, the end-to-end delay of p_f^j is given by

$$d_f^j \leq GRC^K(p_f^j) + \alpha^K - A^1(p_f^j) \quad (23)$$

where $\alpha^K = \beta^K + \tau^{K,K+1}$ and $\tau^{K,K+1}$ is the propagation delay between server K and the destination.

Observe that $GRC^K(p_f^j)$ depends on $A^K(p_f^j)$, which in turn depends on $GRC^{K-1}(p_f^j)$. Applying this argument recursively $GRC^K(p_f^j)$ can be related to $GRC^1(p_f^j)$. Since $GRC^1(p_f^j)$ is completely determined by the arrival characteristics of the source and the rate associated with the flow, the end-to-end delay can be determined if source specification is known. In what follows, we first relate $GRC^{i+1}(p_f^j)$ to $GRC^i(p_f^j)$ and then relate the end-to-end delay to $GRC^1(p_f^j)$.

Lemma 1. *If the scheduling algorithm at server i belongs to GR for flow f , then*

$$GRC^{i+1}(p_f^j) \leq GRC^i(p_f^j) + \max_{k \in [1 \dots j]} \frac{l_f^k}{r_f} + \alpha^i \quad j \geq 1 \quad (24)$$

where $\alpha^i = \beta^i + \tau^{i,i+1}$ and $\tau^{i,i+1}$ is the propagation delay between servers i and $i + 1$.

Proof. The proof is by induction on j .

Base Case: $j = 1$

$$GRC^{i+1}(p_f^1) = A^{i+1}(p_f^1) + \frac{l_f^1}{r_f} \quad (25)$$

Since scheduling algorithm at server i belongs to GR for flow f , $A^{i+1}(p_f^1) \leq GRC^i(p_f^1) + \alpha^i$. Hence,

$$GRC^{i+1}(p_f^1) \leq GRC^i(p_f^1) + \frac{l_f^1}{r_f} + \alpha^i \quad (26)$$

$$\leq GRC^i(p_f^1) + \max_{k \in [1 \dots 1]} \frac{l_f^k}{r_f} + \alpha^i \quad (27)$$

Therefore (24) holds for $j = 1$

Induction hypothesis: Assume (24) holds for $1 \leq j \leq m$.

Induction: We need to show (24) holds for $1 \leq j \leq m + 1$.

$$\begin{aligned} GRC^{i+1}(p_f^{m+1}) &= \max\{A^{i+1}(p_f^{m+1}), GRC^{i+1}(p_f^m)\} + \frac{l_f^{m+1}}{r_f} \end{aligned} \quad (28)$$

Since scheduling algorithm at server i belongs to GR for flow f , $A^{i+1}(p_f^{m+1}) \leq GRC^i(p_f^{m+1}) + \alpha^i$. Hence,

$$\begin{aligned} GRC^{i+1}(p_f^{m+1}) &\leq \max\{GRC^i(p_f^{m+1}) \\ &+ \alpha^i, GRC^{i+1}(p_f^m)\} + \frac{l_f^{m+1}}{r_f} \end{aligned} \quad (29)$$

Thus, there are two cases to consider:

1. If $GRC^i(p_f^{m+1}) + \alpha^i > GRC^{i+1}(p_f^m)$, then from (29) we get

$$\begin{aligned} GRC^{i+1}(p_f^{m+1}) &\leq GRC^i(p_f^{m+1}) + \frac{l_f^{m+1}}{r_f} + \alpha^i \end{aligned} \quad (30)$$

Since $\max_{k \in [1 \dots m+1]} \frac{l_f^k}{r_f} \geq \frac{l_f^{m+1}}{r_f}$

$$\begin{aligned} GRC^{i+1}(p_f^{m+1}) &\leq GRC^i(p_f^{m+1}) + \max_{k \in [1 \dots m+1]} \frac{l_f^k}{r_f} + \alpha^i \end{aligned} \quad (31)$$

2. If $GRC^i(p_f^{m+1}) + \alpha^i \leq GRC^{i+1}(p_f^m)$, then from (29) we get

$$GRC^{i+1}(p_f^{m+1}) \leq GRC^{i+1}(p_f^m) + \frac{l_f^{m+1}}{r_f} \quad (32)$$

Using the induction hypothesis we get,

$$\begin{aligned} GRC^{i+1}(p_f^{m+1}) &\leq GRC^i(p_f^m) + \max_{k \in [1 \dots m]} \frac{l_f^k}{r_f} + \frac{l_f^{m+1}}{r_f} + \alpha^i \end{aligned} \quad (33)$$

Since $GRC^i(p_f^{m+1}) \geq GR C^i(p_f^m) + \frac{l_f^{m+1}}{r_f}$ and $\max_{k \in [1 \dots m+1]} \frac{l_f^k}{r_f} \geq \frac{l_f^{m+1}}{r_f}$,

$$GRC^{i+1}(p_f^{m+1}) \leq GR C^i(p_f^{m+1}) + \max_{k \in [1 \dots m+1]} \frac{l_f^k}{r_f} + \alpha^i \quad (34)$$

From (31), (34) and the induction hypothesis, we conclude that (24) holds for $1 \leq j \leq m+1$. Hence, the lemma follows.

Theorem 2. *If the scheduling algorithm at each of the servers on the path of a flow belongs to GR for flow f , then the end-to-end delay of packet p_f^j , denoted by d_f^j , is given by*

$$d_f^j \leq GR C^1(p_f^j) - A^1(p_f^j) + (K-1) \max_{n \in [1 \dots j]} \frac{l_f^n}{\gamma_f} + \sum_{n=1}^{n=K} \alpha^n \quad (35)$$

where $\alpha^n = \beta^n + \tau^{n,n+1}$ and K is the number of servers on the path of the flow.

Proof. Since server K guarantees that packet p_f^j will be transmitted by time $GR C^K(p_f^j) + \beta^K$ and the packet arrives at the first node at time $A^1(p_f^j)$,

$$d_f^j \leq GR C^K(p_f^j) + \alpha^K - A^1(p_f^j) \quad (36)$$

Since the scheduling algorithm at each server on the path of the flow belongs to GR, by repeated application of Lemma 1, we conclude that:

$$GR C^K(p_f^j) \leq GR C^1(p_f^j) + (K-1) \max_{n \in [1 \dots j]} \frac{l_f^n}{\gamma_f} + \sum_{n=1}^{n=K-1} \alpha^n \quad (37)$$

The theorem follows from (36) and (37).

Notice that $\alpha^n = \beta^n + \tau^{n,n+1}$. Hence, it is completely characterized by the scheduling algorithm and the propagation delay in the network. The remaining terms in (35) depend on the source traffic specifications, and are evaluated in the next section.

3.2 Source traffic specifications

A source traffic is specified by characterizing the number of bits that arrive at the network over an interval of time. Let $AP_f^i(t_1, t_2)$ be a function that denotes the flow f bits that arrive in the interval $[t_1, t_2]$ at server i . The bits of a packet are considered to have arrived only after the arrival of the last bit of the packet. Hence, the arrival function consists of an impulse at each packet arrival instance and is right continuous. Also $AP_f^i(t, t)$ is the length of the packet that arrives at time instant t .

As may be evident from Theorem 2, to determine an upper bound on end-to-end packet delays, we must relate the the Guaranteed Rate clock value of a packet to its arrival time at the first server. To achieve this objective, we define set S_f^i for flow f at server i as follows:

$$S_f^i = \{n \mid n > 0 \wedge GR C^i(p_f^{n-1}) \leq A^i(p_f^n)\}$$

From the definitions of GR clock and set S_f^i , it can be shown that for each packet p_f^j

$$GR C^i(p_f^j) = A^i(p_f^k) + \sum_{n=0}^{n=j-k} \frac{l_f^{k+n}}{r_f} \quad (38)$$

where $k \leq j$ is the largest integer belonging to set S_f^i . Since $\sum_{n=0}^{n=j-k} l_f^{k+n} = AP_f^i(A^i(p_f^k), A^i(p_f^j))$, we get

$$GR C^i(p_f^j) = A^i(p_f^k) + \frac{AP_f^i(A^i(p_f^k), A^i(p_f^j))}{r_f} \quad (39)$$

Using Theorem 2 and (39), we will now determine deterministic end-to-end delay bound for Leaky Bucket sources (a deterministic source characterization) and upper bound on the tail distribution of the delay for EBB sources (a stochastic characterization).

3.2.1 Leaky bucket

Leaky Bucket is a source traffic specification that bounds the maximum deviation from the average rate. Specifically, a flow f conforms to Leaky Bucket [13] with burst size σ_f and average rate r_f if

$$AP_f(t_1, t_2) \leq \sigma_f + r_f(t_2 - t_1) \quad t_1 \leq t_2, t_1 \geq 0 \quad (40)$$

Theorem 3. *If flow f conforms to a Leaky Bucket with parameters (σ_f, r_f) and the scheduling algorithm at each of the server on the path of a flow belongs to GR for the flow, then the end-to-end delay of packet p_f^j , denoted by d_f^j , is given by*

$$d_f^j \leq \frac{\sigma_f + (K-1) \max_{n \in [1 \dots j]} l_f^n}{r_f} + \sum_{n=1}^{n=K} \alpha^n \quad (41)$$

where $\alpha^n = \beta^n + \tau^{n,n+1}$ and K is the number of servers on the path of the flow.

Proof. Let $k \leq j$ be largest integer belonging to set S_f^1 . Clearly, such a k must exist. Since flow f conforms to the Leaky Bucket specification, we get

$$AP_f^1(A^1(p_f^k), A^1(p_f^j)) \leq \sigma_f + r_f(A^1(p_f^j) - A^1(p_f^k)) \quad (42)$$

Hence, from (39) we get,

$$GR C^1(p_f^j) \leq \frac{\sigma_f}{r_f} + A^1(p_f^j) \quad (43)$$

which implies

$$GR C^1(p_f^j) - A^1(p_f^j) \leq \frac{\sigma_f}{r_f} \quad (44)$$

Theorem 3 follows from (44) and Theorem 2.

If the scheduling algorithm used at each of the servers is either Virtual Clock or PGPS, then the bound derived in Theorem 3 is tighter than the previously known bounds for rate-proportional processor sharing (RPPS) rate assignment of PGPS networks when the sources conform to Leaky

Bucket [12, 13]. Specifically, the bound in [12] for Leaky Bucket sources is:

$$d_f^j = \frac{\sigma + 2 * (K - 1) * l_f^{\max}}{\rho} + \sum_{n=1}^{n=K} \alpha^n \quad (45)$$

Clearly, the factor 2 in (45) makes the bound significantly loose as compared to the bound in (41). Since (41) improves upon the only term that depends on the network, the improvement is significant.

3.2.2 Exponentially Bounded Burstiness

Exponentially bounded burstiness (EBB) process characterization has been proposed as a statistical relaxation of Leaky Bucket. A flow conforms to EBB if the probability of deviation of a source from the average rate decreases exponentially. Specifically, a flow f conforms to an EBB [16] process with parameters $(r_f, \Lambda_f, \gamma_f)$, if

$$\begin{aligned} Pr(AP_f^i(t_1, t_2) \geq r_f(t_2 - t_1) + x) \\ \leq \Lambda_f e^{-\gamma_f x} \quad x \geq 0 \quad t_1 \leq t_2 \end{aligned} \quad (46)$$

Theorem 4. *If flow f conforms to EBB with parameters $(r_f, \Lambda_f, \gamma_f)$ and the scheduling algorithm at each of the servers on the path of a flow belong to GR for the flow, then the end-to-end delay of packet p_f^j denoted by d_f^j is given by*

$$\begin{aligned} Pr \left(d_f^j \geq x + (K - 1) \frac{\max_{n \in \{1 \dots j\}} l_f^n}{r_f} + \sum_{n=1}^{n=K} \alpha^n \right) \\ \leq \Lambda_f e^{-\gamma_f x r_f} \quad x \geq 0 \end{aligned} \quad (47)$$

where $\alpha^n = \beta^n + \tau^{n, n+1}$ and K is the number of servers on the path of the flow.

Proof. Let $k \leq j$ be the largest integer belonging to S_f^1 . Clearly, such a k must exist. By the definition of EBB process, we have

$$\begin{aligned} Pr \left(AP_f^1(A^i(p_f^k), A^1(p_f^j)) \geq r_f(A^1(p_f^j) - A^1(p_f^k)) + y \right) \\ \leq \Lambda_f e^{-\gamma_f y} \end{aligned} \quad (48)$$

which implies

$$\begin{aligned} Pr \left(\frac{AP_f^1(A(p_f^k), A^1(p_f^j))}{r_f} \geq A^1(p_f^j) - A^1(p_f^k) + \frac{y}{r_f} \right) \\ \leq \Lambda_f e^{-\gamma_f y} \end{aligned} \quad (49)$$

From (39) we have

$$\begin{aligned} Pr \left(GRC^1(p_f^j) - A^1(p_f^j) \geq \frac{y}{r_f} \right) \\ \leq \Lambda_f e^{-\gamma_f y} \end{aligned} \quad (50)$$

The theorem follows from (50) and Theorem 2.

3.3 Discussion

Theorem 2 demonstrates that a heterogeneous network of servers, each of which employs a scheduling algorithm in

the GR class, can provide an upper bound on end-to-end packet delays for heterogeneous sources. Additionally, the salient features of the class of GR scheduling algorithms and the method for determining end-to-end delay bounds include:

- An end-to-end delay bound can be determined for any source traffic specification for which a bound on difference between the Guaranteed Rate clock and the arrival time of a packet at the first server can be determined. Since this difference is only a function of the source traffic characteristics and the rate associated with the flow, it is simple to determine an end-to-end delay bound. Moreover, this difference can be interpreted as delay experienced by a packet at a single queue server with capacity r_f . Hence, queuing analysis can also be used to determine a bound on the tail distribution of delays experienced by packets for a conventional stochastic process characterization.
- A source can determine end-to-end delays without specifying the shape of the traffic to the network by keeping track of Guaranteed Rate clock values associated with its flow at the first server (the reserved rate and other constants, i.e., $\alpha^{i, i+1}$, required for this can be determined at the flow setup time). Such a capability is important, as the source may not have a good characterization of the traffic or the characterization may not be known a priori. Moreover, even if the characterization is known, it may not conform to the set of characterizations supported by the network. Finally, this capability allows the source to decide when to renegotiate the reserved rate.

The method derives its simplicity from the constraint placed on the rate assignments in the scheduling algorithms. Note that the constraint we have placed on rate assignments for PGPS so that it belongs to GR is similar to RPPS [13]. Such constraints allow each flow to be considered in isolation and hence facilitate the determination of delay bounds for each flow independently. If this constraint is relaxed to derive delay bounds for more general rate assignments, assumptions about the behavior of the other sources in the network may have to be made². In such scenarios, however, the QoS guarantees provided to each flow may be conditional since sources may be greedy and not conform to the specifications [14]. Moreover, if hardware traffic enforcement mechanisms are used, the probability of failure of one of the numerous pieces of enforcement hardware may not be negligible, thereby weakening the guarantees. Hence, we believe that the constraints on rate assignments described above are not a limitation, but in fact are highly desirable.

4 Conclusions

In this paper, we have defined a class of GR scheduling algorithms. The GR class includes Virtual Clock, Self-Clocked

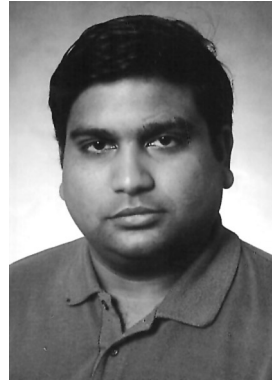
² General rate assignments such as consistent relative session treatment (CRST) for PGPS and SCFQ networks were designed to separate rate from delay allocation. PGPS and SCFQ with CRST rate assignment do not belong to GR. However, the functionality of separation of delay from throughput can be achieved by Delay EDD [18] (for flows with constant packet length) and Leave-in-Time [7] (for flows with variable packet length), both of which can be shown to belong to GR.

Fair Queuing and Packet-by-Packet Generalized Processor Sharing. For networks that employ scheduling algorithms belonging to GR, we have presented a method for determining an upper bound on end-to-end packet delays. The method facilitates determining end-to-end delay for a variety of sources. We have illustrated the method by determining end-to-end delay for sources conforming to Leaky Bucket and Exponentially Bounded Burstiness. The delay bounds that we have derived are parametrized by the scheduling algorithms which, when instantiated with Virtual Clock and SCFQ scheduling algorithms, lead to many new results.

Acknowledgements. This research was supported in part by IBM Graduate Fellowship, the National Science Foundation (Research Initiation Award CCR-9409666), National Science Foundation Grant No. NCR-9004464, NASA, Mitsubishi Electric Research Laboratories (MERL), and Sun Microsystems Inc.

References

1. Clark DD, Shenker S, Zhang L (1992) Supporting real-time applications in an integrated services packet network. In: Proceedings of ACM SIGCOMM, pp 14–26
2. Cruz RL (1991) A calculus for network delay, Part I: network elements in isolation. *IEEE Trans Inf Theory* 37:114–131
3. Cruz RL (1991) A calculus for network delay, Part II: network analysis. *IEEE Trans Inf Theory*, 37:132–141
4. Davin J, Heybey A (1990) A simulation study of fair queuing and policy enforcement. *Comput Commun Rev* 20(5):23–29
5. Demers A, Keshav S, Shenker S (1989) Analysis and simulation of a fair queuing algorithm. In: Proceedings of ACM SIGCOMM, pp 1–12
6. Ferrari D, Verma DC (1990) A scheme for real-time channel establishment in wide-area networks. *IEEE J Sel Areas Commun* 8:368–379
7. Figuera N, Pasquale J (1995) Leave-in-time: a new service discipline for real-time communication in a packet-switching data network. In: *ACM SIGCOMM 95*, pp 207–218
8. Golestani SJ (1991) A framing strategy for congestion management. *IEEE J Sel Areas Commun* pp 1064–1077
9. Golestani SJ (1994) A self-clocked fair queuing scheme for high-speed applications. In: Proceedings of INFOCOM '94
10. Kalmanek CR, Kanakia H, Keshav S (1990) Rate-controlled servers for very high-speed networks. In: Proceedings of IEEE GLOBECOM '90, San Diego, CA, pp 300.3.1–300.3.9
11. Lam SS, Xie GG (1995) Burst scheduling: architecture and algorithm for switching packet video. In: Proceedings of INFOCOM '95
12. Parekh AK, Gallager RG (1994) A generalized processor-sharing approach to flow control in integrated services networks: the multiple node case. *IEEE/ACM Trans Networking* 2(2):137–150
13. Parekh AK (1992) A generalized processor-sharing approach to flow control in integrated services networks. PhD thesis, Department of Electrical Engineering and Computer Science, MIT, Cambridge, Mass.
14. Shenker S (1994) Making Greed Work in Networks: A Game-Theoretic Analysis of Switch Service Disciplines. In: Proceedings of ACM SIGCOMM '94, pp 47–57
15. Xie GG, Lam SS (1994) Delay Guarantee of Virtual Clock Server. Technical Report TR-94-24, Dept. of Computer Sciences, UT-Austin, TX
16. Yaron O, Sidi M (1994) Generalized processor-sharing networks with exponentially bounded burstiness arrivals. In: Proceedings of INFOCOM '94
17. Zhang H, Ferrari D (1993) Rate-controlled static priority queuing. In: Proceedings of INFOCOM '93, volume 2, pp 227–236
18. Zhang H, Keshav S (1991) Comparison of rate-based service disciplines. In: Proceedings of ACM SIGCOMM, pp 113–121
19. Zhang L (1990) Virtual clock: a new traffic control algorithm for packet-switching networks. In: Proceedings of ACM SIGCOMM '90, pp 19–29



PAWAN GOYAL received his B.Tech. in Computer Sciences from the Indian Institute of Technology, Kanpur, in 1992, and is currently a doctoral candidate at the Department of Computer Sciences at the University of Texas at Austin. His research interests are in network, file systems and operating system services for multimedia applications. He has received several awards for academic excellence including the MCD Fellowship (awarded by the University of Texas at Austin), the IBM Doctoral Fellowship, and the Core Year Award from the Indian Institute of Technology, Kanpur.



SIMON S. LAM received the BSEE degree, with distinction, from Washington State University in 1969, and the M.S. and Ph.D. degrees in engineering from the University of California at Los Angeles in 1970 and 1974, respectively. From 1971 to 1974, he was a post-graduate research engineer at the ARPA Network Measurement Center (UCLA). From 1974 to 1977, he was a research staff member at the IBM T.J. Watson Research Center, Yorktown Heights, New York. Since 1977, he has been on the faculty of the University of Texas at Austin, where he is a Professor of Computer Sciences. He holds two anonymously endowed professorships, and served as department chair from 1992 to 1994. His research interests are in network protocol design, performance analysis, formal verification, network security, and multimedia. He was elected an IEEE Fellow in 1985, and served on the editorial boards of *IEEE Transactions on Software Engineering*, *IEEE Transactions on Communications*, *Performance Evaluation*, and *Proceedings of the IEEE*. He presently serves as editor-in-chief of *IEEE/ACM Transactions on Networking*.



HARRICK M. VIN received his Ph.D. in Computer Science from the University of California at San Diego in 1993. He is currently an Assistant Professor of Computer Sciences, and the Director of the Distributed Multimedia Computing Laboratory at the University of Texas at Austin. His research interests are in the areas of multimedia systems, high-speed networking, mobile computing, and large-scale distributed systems. Over the past 5 years, he has coauthored more than 55 papers in leading journals and conferences in the area of multimedia systems. He is a member of the editorial board of *IEEE Multimedia*, the conference and program co-chair for the Multimedia Computing and Networking 1996 (MMCN96) and Multimedia Computing and Networking 1997 (MMCN97), vice-chair in the area of Distributed Multimedia Systems for the 17th International Conference on Distributed Computing Systems (ICDCS), and a program committee member of several conferences, including ACM Multimedia'96, ACM Multimedia'95, and Second International Conference on Distributed Multimedia Systems and Applications. During his career, he has been a recipient of several awards including the National Science Foundation CAREER award, the IBM Faculty Development Award, the National Science Foundation Research Initiation Award, the IBM Doctoral Fellowship, the NCR Innovation Award, and the San Diego Supercomputer Center Creative Computing Award.