

A GENERALIZED PROCESSOR SHARING APPROACH TO FLOW CONTROL IN INTEGRATED SERVICES NETWORKS—THE MULTIPLE NODE CASE

Abhay K. Parekh
T. J. Watson Research Center
International Business Machines
Yorktown Heights NY 01598

Robert G. Gallager
Department of EECS
Massachusetts Institute of Technology
Cambridge MA 02139

Abstract

Worst-case bounds on delay and backlog are derived for leaky bucket constrained sessions in arbitrary topology networks of Generalized Processor Sharing (GPS) [8] servers. When only a subset of the sessions are leaky bucket constrained, we give succinct per-session bounds that are independent of the behavior of the other sessions and also of the network topology. However, these bounds are only shown to hold for each session that is guaranteed a backlog clearing rate that exceeds the token arrival rate of its leaky bucket.

When all of the sessions are leaky bucket constrained, a much larger class of networks called Consistent Relative Session Treatment (CRST) networks is analyzed. The session i route is treated as a whole, yielding tighter bounds than those that result from adding the worst-case delays (backlogs) at each of the servers in the route. The bounds on delay and backlog for each session are efficiently computed from a universal service curve, and it is shown that these bounds are achieved by “staggered” greedy regimes when an independent sessions relaxation holds. Propagation delay is also incorporated into the model. Finally, the analysis of arbitrary topology GPS networks is related Packet GPS networks (PGPS).

1 Introduction

The problem of providing performance guarantees to the diverse users of an integrated services network is central to supporting real-time services such as voice and video. This problem is especially difficult in the presence of congestion, when it is important to use the link bandwidth efficiently. In [8] we proposed the combination of leaky bucket admission control and a work-conserving packet service discipline at the nodes of the network, to accommodate the delay and throughput requirements of a wide range of co-existing sessions. The service discipline is based on Generalized Processor Sharing (GPS) and was first suggested in [3] in the context of managing congestion at gateway nodes. The emphasis there was on treating the users equally; our focus is on the inherent flexibility of the mecha-

nism, and on providing good network-wide per-session bounds on worst-case delay and backlog.

In [8] we analyzed single node systems; here we extend this analysis to arbitrary topology networks of GPS servers. The GPS results are then related to networks in which the nodes follow a packet-based service discipline, Packet GPS (PGPS) discussed extensively in [8]. Due to considerations of space, many proofs are omitted in this paper—the interested reader is referred to [7].

A GPS server that serves N sessions on a link is characterized by N positive real numbers, $\phi_1, \phi_2, \dots, \phi_N$. These numbers denote the relative amount of service to each session in the sense that if $S_i(\tau, t)$ is defined as the amount of session i traffic served during an interval $[\tau, t]$, then

$$\frac{S_i(\tau, t)}{S_j(\tau, t)} \geq \frac{\phi_i}{\phi_j}, \quad j = 1, 2, \dots, N \quad (1)$$

for any session i that is backlogged in the interval $[\tau, t]$. Thus (1) is satisfied with equality for two sessions i and j that are both backlogged during the interval $[\tau, t]$.

Note from (1) that whenever session i is backlogged it is guaranteed a service rate of

$$g_i = \frac{\phi_i}{\sum_{j=1}^N \phi_j} r, \quad (2)$$

where r is the rate of the link. This rate is called the *session i backlog clearing rate* since a session i backlog of size q is served in at most $\frac{q}{g_i}$ time units.

We assume a virtual circuit, connection-based packet network, and analyze the performance of leaky bucket constrained sessions. The session i leaky bucket is characterized by a token bucket of size σ_i and a token arrival rate of ρ_i . The amount of session i traffic entering the network during any interval $(\tau, t]$ is defined to be $A_i(\tau, t)$; if session i is leaky bucket constrained, then

$$A_i(\tau, t) \leq \sigma_i + \rho_i(t - \tau), \quad \forall t \geq \tau \geq 0. \quad (3)$$

As in [8], we say that A_i conforms to (σ_i, ρ_i) , or $A_i \sim (\sigma_i, \rho_i)$. For details on how to accommodate peak rate constraints as well, see [6]. The constraint (3) is identical to the one suggested by Cruz [1].

The main question we address in this paper is the following: Given a network with the values of the server parameters fixed and a set of leaky bucket constrained sessions, what is the worst-case session delay and backlog for each of the sessions in this set?

In Section 2 we set up our model of the network and specify notation. Then the notions of network backlog and delay are discussed and graphically interpreted. Section 4 contains succinct per-session bounds for the leaky bucket constrained sessions of a network, which are independent of the topology and of the behavior of other sessions. Next, we treat the case when all of the sessions are leaky bucket constrained. An important tool for the analysis, the All-Greedy bound, is presented in Section 6. In Section 7, an algorithm is derived that enables a characterization of internal traffic in terms of burstiness, average and peak rates for a broad class of server allocations called Consistent Relative Session Treatment (CRST) assignments. This class of assignments is flexible enough to accommodate a wide variety of session delay constraints. In Section 8, we show that worst-case session delay and backlog can be bounded from an easily computable universal service curve. This is accomplished even though *different* worst-case regimes may maximize delay and backlog for a given session. The bounds are shown to be tight under an independent relaxation assumption, when the traffic follows *staggered greedy* regimes. Propagation delay is incorporated in Section 9, and our results are extended to PGPS networks in Section 10. Conclusions are in Section 11.

Note that all of our bounds can be applied to networks of arbitrary topology.

2 The Network Model

The network is modeled as a directed graph in which nodes represent switches and arcs represent links. A route is a path in the graph, and the path taken by session i is defined as $P(i)$. Let $P(i, k)$ be the k^{th} node in $P(i)$, and K_i be the total number of nodes in $P(i)$. The rate of the server at node m is r^m .

The amount of session i traffic that enters the network in the interval $(\tau, t]$ is given by $A_i(\tau, t)$. Let $S_i^{(k)}(\tau, t)$, $k = 1, \dots, K_i$, be the amount of session i traffic served by node $P(i, k)$ in the interval $[\tau, t]$. Thus, $S_i^{(K_i)}$ describes the traffic that leaves the network. We characterize the service function by “pseudo” leaky bucket parameters $\sigma_i^{(k)}$ and ρ_i so that

$$S_i^{(k)}(\tau, t) \leq \sigma_i^{(k)} + \rho_i(t - \tau), \quad \forall t \geq \tau \geq 0, \quad (4)$$

i.e., $S_i^{(k)} \sim (\sigma_i^{(k)}, \rho_i)$.

Often, we will analyze what happens at a particular server, m . In this case the notation described above becomes overly cumbersome. Define $I(m)$ to be the

set of sessions that are served by server m . For every session $i \in I(m)$, let the arrival function into that node be described by $A_i^m \sim (\sigma_i^m, \rho_i)$ and the departure function be described by $S_i^m \sim (\sigma_i^{m, \text{out}}, \rho_i)$. For example, at server 0 in Figure 1: $A_0^0 = A_0$, $A_2^0 = S_2^{(2)}$, and $A_3^0 = S_3^{(1)}$. Thus when $k = P(i, j)$ for a particular session, i , the functions $S_i^{(j)}$ and S_i^k are identical. The

Figure 1: A four server network. The demultiplexer works instantaneously.

rate of the link associated with server m is denoted by r^m . The value of ϕ_i at m is denoted by ϕ_i^m for all $i \in I(m)$. Finally, let g_i^m be the session i backlog clearing rate from (2) at node m , i.e.,

$$g_i^m = \frac{\phi_i^m}{\sum_{j \in I(m)} \phi_j^m} r^m. \quad (5)$$

3 Network Delay, Backlog and Stability

In this section we extend the notions of session i delay and backlog introduced in [8] to the multiserver case. Given a set of arrival functions for every session in the network, define $Q_i^{(k)}(t)$ to be the session i backlog at node $P(i, k)$ at time t . Similarly, let $Q_i^m(t)$ be the session i backlog at node $m \in P(i)$. Thus, if $m = P(i, k)$, then

$$Q_i^{(k)}(t) = Q_i^m(t) = A_i^m(0, t) - S_i^m(0, t). \quad (6)$$

Define the total session i backlog at time t to be

$$Q_i(t) = \sum_{k=1}^{K_i} Q_i^{(k)}(t). \quad (7)$$

Thus, $Q_i(t)$ is the amount of session i traffic buffered in the network at time t . By assumption,

$$Q_i(t) = 0, \quad \forall t \leq 0$$

for every session i . Also, let $D_i(t)$ be the time spent in the network by a session i bit that arrives at time t . Figure 2 shows how to represent the notions of backlog and delay graphically. The top figure shows how session i traffic progresses through the nodes of its route. Notice that the arrival function to node 2 is the session i service function of node 1.

The bottom figure shows how the backlog and delay can be measured and illustrates the definitions of Section 3. We see that $D_i(\tau)$ is the horizontal distance

Figure 2: An Example of Session i flow when $K_i = 2$.

between the curves $A_i(0, t)$ and $S_i^{(K_i)}(0, t)$ at the ordinate value of $A_i(0, \tau)$. Clearly, $D_i(\tau)$ depends on the arrival functions A_1, \dots, A_N , where N is the total number of sessions in the network). We are interested in computing the maximum delay over all time, and over all arrival functions that are consistent with (3). Let D_i^* be the maximum delay for session i . Then

$$D_i^* = \max_{(A_1, \dots, A_N)} \max_{\tau \geq 0} D_i(\tau).$$

The backlogs at every node in $P(i)$ can be determined from Figure 2 as shown.

Define the maximum backlog for session i , Q_i^* :

$$Q_i^* = \max_{(A_1, \dots, A_N)} \max_{\tau \geq 0} Q_i(\tau).$$

Note that in Figure 2 A_i contains an impulse at time a ; As in [8], we adopt the convention that the arrival functions are continuous from the left, so that $A_i(0, a) = \alpha$ and $A_i(0, a^+) = \hat{\alpha}$.

Define the utilization of server m to be

$$u^m = \frac{\sum_{j \in I(m)} \rho_j}{r^m}. \quad (8)$$

A network is defined to be *stable* if $D_i^* < \infty$ for all sessions i . In most of our analysis we will show stability

under the assumption that $u^m < 1$ at every server m . Allowing utilizations of greater than 1 would permit backlogs and delays to build up unboundedly, and we have shown elsewhere ([6]) that permitting $u^m = 1$ at each server m can result in problems as well.

The minimum session i backlog clearing rate along its route is

$$g_i = \min_{m \in P(i)} g_i^m. \quad (9)$$

When $g_i > \rho_i$ we define session i to be *locally stable*. Note that if $\phi_i^m = \rho_i$ and $u^m < 1$ for all i and $m \in P(i)$, then each session i is locally stable.

Finally, the definitions of system and session busy periods given in [8] for a single node are extended to the multiple node case. A network system (session i) busy period is defined to be the maximal interval B (B_i) such that for every $\tau \in B$ ($\tau \in B_i$), there is at least one server in the network that is in a system (session i) busy period at time τ .

4 Bounds for Locally Stable Sessions

While every route in a data network is acyclic, the union of several routes may result in cycles being induced in the network topology. The presence of these cycles can complicate the analysis of delay considerably, but more importantly, it can lead to feedback effects that drive the system towards instability. This phenomenon has been noticed by researchers from fields as diverse as manufacturing systems [9, 5], communication systems [2] and VLSI circuit simulation [4]. Consider the four node example in Figure 1 (which is identical to Example 2 of Cruz [2]). Suppose the service discipline is FCFS. As an illustration of virtual feedback, notice that $S_0^{(1)}$ depends on the traffic from sessions 2, 3, 4, but the form of this traffic is not independent of $S_0^{(1)}$.

In this section we will show that for a locally stable session, i , these virtual feedback effects are completely absent even when the other sessions are not leaky bucket constrained. For notational convenience let $P(i) = (1, 2, \dots, K_i)$. The following useful Lemma is straightforward and stated without proof—to see that it is true, recall that we are ignoring propagation delays:

Lemma 1 *For every interval $[\tau, t]$ that is contained in a single session i network busy period:*

$$S_i^{(K_i)}(\tau, t) \geq g_i (t - \tau).$$

The Lemma leads us to the main result of this section:

Theorem 1 *If $g_i \geq \rho_i$ for session i :*

$$Q_i^* \leq \sigma_i,$$

$$D_i^* \leq \frac{\sigma_i}{g_i}.$$

Note that the delay bound in Theorem 1 is independent of the topology of the network and number of links in the route taken by the session. Also, it is independent of the σ_j , $j \neq i$.

Proof. Suppose Q_i^* is achieved at time t , and let τ be the first time before t when there are no session i bits backlogged in the network. Then by Lemma 1, $S_i^{(K_i)}(\tau, t) \geq \rho_i(t - \tau)$. Consequently,

$$Q_i^* \leq (\sigma_i + \rho_i(t - \tau)) - \rho_i(t - \tau) = \sigma_i.$$

An arriving session i bit will be served after at most Q_i^* session i bits have been served. Using Lemma 1 again, these backlogged bits are served at a rate of at least g_i . Therefore:

$$D_i^* \leq \frac{Q_i^*}{g_i} \leq \frac{\sigma_i}{g_i}.$$

□

Notice that the bounds are independent of K_i , the number of hops in the session i route. The naive bound on delay arrived at by adding the worst-case delays at each node is $D_i^* \leq \sigma_i \sum_{m=1}^{K_i} \frac{1}{g_i^m}$, illustrating the fact that much better bounds result from analyzing the session i route as a whole. When all of the sessions are leaky bucket constrained and $\phi_i^m = \rho_i$ ($u_i^m < 1$) at all m and $i \in I(m)$:

$$Q_i^* \leq \sigma_i, \quad (10)$$

and

$$D_i^* \leq \frac{\sigma_i}{\rho_i}. \quad (11)$$

However, note that given a locally stable session i , the result of Theorem 1 is valid for any GPS assignment for the other sessions. In fact, the other sessions need not be leaky bucket constrained, nor need the system be stable.

5 The Importance of Sessions that are not Locally Stable

When **all** of the sessions are leaky bucket constrained, it is possible to guarantee finite delay even for the sessions that are not locally stable. This is because GPS is work conserving and the token arrival rates are assigned such that $\sum_{j \in I(m)} \rho_j < r^m$ at all nodes m . Thus we may allow g_i to be less than ρ_i for sessions that are not delay sensitive, and much greater than ρ_i for delay sensitive sessions.

To see why such assignments are important, consider the following example illustrated in Figure 3: There are two sessions in the network, and $P(i) =$

Figure 3: Giving delay sensitive, steady sessions large values of ϕ .

(1, 2) for $i = 1, 2$. Session 1 is more steady than Session 2. Although the backlog clearing rate for the Session 2 is infinitesimal when $\phi_1^m \gg \phi_2^m$, $m = 1, 2$, the session 2 delay is not increased significantly. Thus by giving session 1 a very large backlog clearing rate we can minimize its delay while degrading session 2 delay only slightly. This shows that even when the backlog clearing rate for session 2 is *much smaller* than ρ_2 at each node, the session does not suffer much in terms of delay.

6 The All-Greedy Bound for a single node

The presence of sessions that are not locally stable complicates our analysis considerably; yet after performing the analysis we will see that the computation of per-session delay and backlog remains efficient and intuitive. There are two steps to providing worst case bounds on delay and backlog: The first consists of characterizing the internal traffic of the network so that at each node, m and $j \in I(m)$ we have σ_j^m such that $A_j^m \sim (\sigma_j^m, \rho_j)$. In the second step, the internal characterization is used to analyze the session i route for delay and backlog.

Central to our analytical technique is the concept of the all-greedy bound: We calculate *upper bounds* on the minimum value $\sigma_i^{m,out}$ such that $S_i^m \sim (\sigma_i^{m,out}, \rho_i)$. These upper bounds will be shown to be quite good for a wide variety of networks. Consider a particular good node m . Suppose that for every $j \in I(m)$, we are given that $A_j^m \sim (\sigma_j^m, \rho_j)$. In [8] it was shown that the worst-case delay and backlog for session i (at node m) is each achieved when all the sessions $j \in I(m)$ are simultaneously greedy from time zero, the beginning of a system busy period. However, if two sessions j and p are both served by the same node, n , just before they contend for node m , then it may

not be possible for both of them to be simultaneously greedy, as is required in the all-greedy regime. Thus, the achievable worst-case delay and backlog at node m may be less (but never more) than that calculated under the all-greedy regime.

In the rest of this paper we will make frequent use of the all-greedy bound, in order to simplify procedures for estimating D_i^* and Q_i^* . The following notation is useful in this regard:

We are given σ_j^m, ρ_j for each $j \in I(m)$, such that $\sum_{j \in I(m)} \rho_j < r^m$. Consider a fictitious system in which no traffic enters node m before time zero, and all the sessions at m are greedy starting at time zero. Denote \hat{A}_i^m as the resulting session i arrival function for all $i \in I(m)$. Also denote \hat{S}_i^m as the service function at node m . Recall from [8], that for $t > 0$, as long as $Q_i^m(t) > 0$, the function $\hat{S}_i^m(0, t)$ is piecewise linear and convex- \cup in t . By using the techniques of [8] we can find the smallest value $\hat{\sigma}_i^{m, out}$ such that $\hat{S}_i^m \sim (\hat{\sigma}_i^{m, out}, \rho_i)$. Proceeding from the discussion above it is not difficult to show that

$$\hat{\sigma}_i^{m, out} \geq \sigma_i^{m, out}. \quad (12)$$

Thus, we may bound the burstiness of S_i^m by $\hat{\sigma}_i^{m, out}$.

7 Non-Acyclic GPS networks under Consistent Relative Session Treatment

We begin by the following useful definition:

Definition. *Session j is said to impede a session i at a node m if*

$$\frac{\phi_i^m}{\phi_j^m} < \frac{\rho_i}{\rho_j}.$$

Note that for any two sessions, i and j , that contend for a node m , either session i impedes session j or vice versa, unless $\frac{\phi_i^m}{\phi_j^m} = \frac{\rho_i}{\rho_j}$, in which case neither session impedes the other.

A Consistent Relative Session Treatment GPS assignment (CRST) is one for which there exists a strict ordering of the sessions such that for any two sessions i, j , if session i is less than session j in the ordering, then session i does not impede session j at any node of the network.

The class of assignments that are CRST is quite broad: For example, consider the special case of a CRST system for which

$$\phi_{ij} = \frac{\phi_i^m}{\phi_j^m}, \quad \forall m \text{ s.t. } i, j \in I(m). \quad (13)$$

Thus, whenever sessions i and j contend for service at a link, they are given the same relative treatment. Note that $\phi_{ij} = \frac{\phi_{ip}}{\phi_{jp}}$, where session p is in $I(i) \cap I(j)$.

Such CRST systems are called Uniform Relative Session Treatment (URST) systems. Note that

- By normalizing the values of the ϕ_i^m 's at each node m , we may equivalently define a URST system to be one in which for every session i , and node m that is on the session i route: $\phi_i = \phi_i^m$.
- Suppose $\phi_i = \rho_i$ for every session i . Then from (9) each session is locally stable. We call this special case of a URST system, Rate Proportional Processor Sharing (RPPS).

We can show that a CRST system is stable if $u^m < 1$ at each node, is stable, and have an algorithm [7] for characterizing the internal traffic for every session in a CRST system. It relies on a result that makes the computation of traffic parameters for a session independent of the traffic parameters of those sessions that do not impede it. The explicit ordering in the definition of CRST allows the algorithm to take advantage of this result. The algorithm is efficient and runs in time $O(\sum_i |P(i)|)$. Some parallelization is also possible and is discussed in [6].

8 Computing Delay and Backlog for Stable Systems with Known Internal Burstiness

Suppose that we are given a stable GPS system in which the sessions are leaky bucket constrained as in (3), i.e., for every session j and node m such that $j \in I(m)$, we are given a value σ_j^m , such that $A_j^m \sim (\sigma_j^m, \rho_j)$. As we discussed in Section 6, worst case delay (backlog) at a single node of the network can be upper bounded by applying the techniques of [8] when the traffic characterization of sessions sharing that node is known. Under the Additive Method due to [2], we add the worst case bounds on delay (backlog) for session i at each of the nodes $m \in P(i)$ considered in isolation. While this approach works for any server discipline for which the single node can be analyzed, it may yield very loose bounds. For example, when applied to an RPPS system (defined in Section 7 we get $D_i^* \leq K_i \frac{\sigma_i}{\rho_i}$, rather than $D_i^* \leq \frac{\sigma_i}{\rho_i}$. The problem, of course, is that we are ignoring strong dependencies among the queueing systems at the nodes in $P(i)$. Figure 4 illustrates the Additive Method: The figures (a) and (b) can be determined independently. The Additive Method yields bounds $Q_i^* \leq Q_i^{1*} + Q_i^{2*}$ and $D_i^* \leq D_i^{1*} + D_i^{2*}$. In order to improve the bounds the session route as a whole is treated as a whole. For notational simplicity we focus on a particular session, i , that follows the route $1, 2, \dots, K$. We will assume that:

1. The sessions $j \in I(m) - \{i\}$ (for $m = 1, 2, \dots, K$) are free to send traffic in any manner as long as

Figure 4: The Additive Method for session i when $P(i) = \{1, 2\}$.

$A_j^m \sim (\sigma_j^m, \rho_j)$. Thus it is appropriate to call the sessions in $I(m) - \{i\}$, the *independent sessions* at node m ($m = 1, 2, \dots, K$).

2. Session i traffic is constrained to flow along its route so that

$$A_i^m = S_i^{m-1} \quad m = 2, 3, \dots, K.$$

Assumptions 1 and 2 are collectively known as the independent sessions relaxation. This is because while the network topology may preclude certain arrival functions of A_j^k that are consistent with (σ_j^k, ρ_j) , these functions are included under the independent sessions relaxation. On the other hand, every arrival function allowable in the network, is allowed under the independent sessions relaxation. Thus, the values of D_i^* and Q_i^* that hold under the independent sessions relaxation, must be upper bounds on the true values of these quantities. The use of all-greedy bounds enables us to compute D_i^* and Q_i^* exactly under the independence relaxation. Figure 5 illustrates the system to be analyzed. In view of our results for the single node

Figure 5: Analyzing the Session i route as a whole, under the Independent Sessions Relaxation.

case, it would be satisfying if maximum delay (and backlog) were achieved when all the sessions of the network are greedy starting at time zero (the beginning of a system busy period). However, this is generally not true. It turns out that what is required

is that the sessions at a particular node j become greedy simultaneously, but only after the sessions at node $j - 1$ become greedy. We call this pattern of arrivals a staggered greedy regime. The instants of time at which the sessions become greedy depend on the session for which maximum delay and backlog is being estimated. We will also find that Q_i^* and D_i^* may not both be achieved for the *same* staggered greedy regime. This important point is illustrated in Figure 6: The curves \hat{S}_i^1 and \hat{S}_i^2 are shown in (a). Note that $\sigma_i^2 = \hat{\sigma}_i^{1,out}$, and so \hat{S}_i^1 and \hat{S}_i^2 cannot be determined independently.

Figure (b) shows two staggered greedy regimes. In the first, the sessions in $I(2) - \{i\}$ become greedy at time t_1 , which yields a maximum backlog of q_i^* at time τ . In the second staggered greedy regime, the sessions at $I(2) = \{i\}$ wait until time \bar{t}_1 to become greedy—this results in a maximum delay of d_i^* for session i at time zero. The possibility of D_i^* and Q_i^* being achieved un-

Figure 6: Two Staggered Greedy Regimes when $P(i) = \{1, 2\}$

der different staggered greedy regimes is discouraging from a practical standpoint, especially if computing either one of these quantities involves solving a complicated optimization problem. It would be much more desirable to have a single function from which both delay and backlog can be bounded. (In the single-node case this curve is just \hat{S}_i , i.e. Lemma 10 of [8].)

In Section 8.1 we describe such a function, which we call the session i universal curve, $U_i(t)$. This curve is constructed without computing any staggered greedy regimes, and both D_i^* and Q_i^* can be determined efficiently and exactly from it (under the independent sessions relaxation). In addition, the staggered greedy regimes that achieve these worst-case values can also be efficiently determined from $U_i(t)$. In Section 8.2, we prove that these worst-case staggered greedy regimes achieve the same bounds on D_i^* and Q_i^* , as computed from $U_i(t)$.

8.1 The Session i Universal Service Curve

For notational simplicity, we will focus on a session i such that $P(i) = (1, 2, \dots, K)$. The functions $\hat{S}_i^1, \dots, \hat{S}_i^K$ (defined in Section 6) can be computed using the internal traffic characterization of Section 7 by using the independent sessions relaxation. Recall that for each node $m = 1, 2, \dots, K$, \hat{S}_i^m is continuous, piece-wise linear and is convex- \cup in the range $[0, t_m^B]$, where t_m^B is the duration of the session i busy period at m under the all-greedy regime. Also $\hat{S}_i^m(0) = 0$. Thus it can be specified (in the range $[0, t_m^B]$) by a list of pairs:

$$(s_1^m, d_1^m), (s_2^m, d_2^m), \dots, (s_{n_m}^m, d_{n_m}^m),$$

where s_j^m is the slope of the j^{th} line segment and d_j^m is its duration. Here

$$s_1^m < s_2^m < \dots < s_{n_m}^m, \quad (14)$$

and

$$\sum_{j=1}^{n_m} d_j^m = t_m^B. \quad (15)$$

We first describe how to construct U_i from $\hat{S}_i^1, \dots, \hat{S}_i^K$, and then define the curve analytically. Finally, we establish the relationship between U_i and the session i departures from the network, $S_i^{(K)}$:

Let E_i^k be the collection of all the pairs (s_j^m, d_j^m) for $m = 1, 2, \dots, k$ —i.e.

$$E_i^k = \bigcup_{m=1}^k \bigcup_{j=1}^{n_m} \{(s_j^m, d_j^m)\}.$$

The session i universal service curve, U_i is defined as:

$$U_i(t) = \min\{G_i^K(t), \hat{A}_i(0, t)\},$$

where the curve G_i^k (for $k = 1, 2, \dots, K$) is a continuous curve constructed from the elements of E_i^k as follows:

1. Set $G_i^k(0) = 0$, Remaining-in-E = E_i^k ; $Glist = \phi$; $u = 0$; $t = 0$.
2. Order the elements of E_i^k in increasing order of slope. Remove from Remaining-in-E an element of smallest slope: $e^{new} = (s^{new}, d^{new})$. Append $Glist$ with e^{new} . If Remaining-in-E is not empty then repeat step 2.
3. G_i^k is a piece-wise linear convex- \cup defined in the range $[0, \sum_{m=1}^k t_m^B]$ by the elements of $Glist^1$. For $t \geq \sum_{m=1}^k t_m^B$ set

$$G_i^k(t) = G_i^k\left(\sum_{m=1}^k t_m^B\right) + \hat{A}_i\left(\sum_{m=1}^k t_m^B, t\right). \quad (16)$$

¹In the same manner as \hat{S}_i^m was specified earlier.

Figure 7 illustrates the construction of U_i for a simple two node example: The two service curves in (a) show \hat{S}_i^1 and \hat{S}_i^2 . In (b) the line segments that make up these curves are concatenated to make a piece-wise linear convex curve that meets \hat{A}_i at time B_K . Note

Figure 7: An example of how U_i is constructed for $K = 2$.

that the line segment with slope s_2^3 is never used in the construction of U_2 , i.e. $T_2 < t_1^B + t_2^B$. The following points are also of interest:

- G_i^k is defined for $k = 1, 2, \dots, K$, but U_i is defined in terms of G_i^K .
- For each m , the relative order of the elements from \hat{S}_i^m is preserved in $Glist$.
- There exists a time B_K such that $U_i(t) = G_i^K(t)$ for $t \leq B_K$ and $U_i(t) = \hat{A}_i^K(t)$ for $t \geq B_K$.

We have shown [7] the following close relationship between S_i^m and G_i^m :

Lemma 2 Consider a given arrival function, A_i , and a given time τ such that $Q_i(\tau) = 0$. Then for each m , $1 \leq m \leq K$, each $t > \tau$:

$$S_i^m(\tau, t) \geq \min_{V \in [\tau, t]} \{A_i(\tau, V) + G_i^m(t - V)\}. \quad (17)$$

In the next section we are going to show that $G_i^m(t)$ is the amount of service given to session i under a specific staggered greedy regime called the (m, t) -staggered greedy regime. Thus Lemma 2 shows that the service to session i is minimized when a such a staggered greedy regime is delayed by an appropriate amount, which is the minimizing value of V . Equation (17) facilitates the following bounds on delay and backlog:

Theorem 2 For every session i :

$$Q_i^* \leq \max_{\tau \geq 0} \{\hat{A}_i(0, \tau) - G_i^K(\tau)\}, \quad (18)$$

and

$$D_i^* \leq \max_{\tau \geq 0} \left\{ \min\{t : G_i^K(t) = \hat{A}_i(0, \tau)\} - \tau \right\}. \quad (19)$$

The inequalities (18) and (19) illustrate the importance of the universal curve. To find the bound on D_i^* compute the maximum horizontal distance between the curves $A_i(0, t)$ and $U_i(t)$ at the ordinate value of $\hat{A}_i(0, t)$. Similarly, Q_i^* is bounded by the maximum vertical distance between the two curves. In the next section, we will show that these bounds are *achieved* for (K, t) -staggered greedy regimes under the independent sessions relaxation.

8.2 The (K, t) -Staggered Greedy Regime

In this section we make clear the relationship between staggered greedy regimes and the session i universal curve U_i . As in the previous sections, we will focus on staggered greedy regimes with respect to a session i and assume that $P(i) = \{1, 2, \dots, K\}$.

Any staggered greedy regime can be characterized by a vector

$$(T_1, \dots, T_K), \quad T_1 \leq T_2 \leq \dots \leq T_K$$

such that all the sessions at node 1 are simultaneously greedy starting at time T_1 , and the independent sessions at node j do not send any traffic in the interval $[T_1, T_j]$, but are simultaneously greedy starting at time T_j . Observe that the first staggered greedy regime in Figure 6(b) can be characterized by $(0, t_1)$ and the second by $(0, \bar{t}_1)$.

A (K, t) -staggered greedy regime, $t \leq B_K$, is the staggered greedy regime characterized by $(0, T_2, \dots, T_K)$ such that

$$\sum_{k=1}^K \hat{S}_i^k(0, T_{k+1} - T_k) = G_i^K(t) \quad (20)$$

where $T_1 = 0$, $T_{K+1} = t$ and $T_{k+1} - T_k \leq t_k^B$ for $k = 1, 2, \dots, K$.

Note that

- Since $t \leq B_K$, $G_i^K(t) = U_i(t)$.
- For each $k = 1, 2, \dots, K - 1$ the staggered greedy regime defined by $(0, T_2, \dots, T_k)$ describes a (k, T_{k+1}) -staggered greedy regime.

The universal service curve can be used to determine T_2, \dots, T_K . This is illustrated in Figure 8 for the simple case of $K = 2$: The top figure the curve U_2 that was constructed from \hat{S}_i^1 and \hat{S}_i^2 . In order to find the $(2, \tau)$ -staggered greedy regime, add the durations of the line segments taken from \hat{S}_i^1 that are in $U_2(t)$, $t \leq \tau$. This sum is T_2 , the time that the independent

Figure 8: Computing a (k, t) -Staggered Greedy Regime when $P(i) = \{1, 2\}$

sessions at node 2 become greedy. This characterizes the staggered greedy regime which is shown in the bottom figure.

Also notice from the figure that in the range $[0, T_2]$, S_i^2 is comprised of the line segments belonging to \hat{S}_i^1 that make up the universal curve in the range $[0, t]$. Also,

$$S_i^2(0, \tau) = U_i(\tau).$$

It turns out that this is true in general:

Theorem 3 For any (K, t) -staggered greedy regime:

$$S_i^{(K)}(0, t) = G_i^K(t).$$

Figure 9 shows how to construct the staggered greedy regimes that maximize backlog and delay. Figure (a)

Figure 9: The staggered greedy regimes that maximize backlog and delay under the independent sessions relaxation.

shows the session i universal curve. Notice that for this curve “backlog” is maximized at time τ_1 and “delay” is maximized at time τ_2 . Figure (b) shows the

two staggered greedy regimes corresponding to these times. Notice that the backlog at time τ_1 in the first regime is exactly equal to the “backlog” at time τ_1 in (a), and similarly the delay at time τ_2 in the second regime is exactly equal to the “delay” at that time in (b).

From Theorems 2 and 3 we have the main theorem of this section:

Theorem 4 *Under the independent sessions relaxation, D_i^* and Q_i^* are each achieved under (K, t) -staggered greedy regimes.*

Now since the values of D_i^* and Q_i^* achieved under the independent sessions relaxation are upper bounds to the actual values of these quantities, we have shown how to find upper bounds on session backlog and delay. Also, since an infinite capacity link can always simulate a finite capacity link, worst case session i backlog and delay calculated under this relaxation must upper bound the values of these quantities for finite capacity links.

9 Propagation Delay

It is easy to incorporate deterministic propagation delays into our network framework. Suppose that every bit transmitted on link (i, j) , incurs a delay of $d_{i,j}$ time units. Then each link acts as a constant delay element, and the characterization of internal traffic (using the method of Section 7) remains the same. A natural modification of the independent sessions relaxation allows us to bound end-to-end delay as well: This leads us to [7]:

$$D_i^* \leq \sum_{m=1}^K d_{m-1,m} + D_i^{*,\text{noprop}},$$

where $D_i^{*,\text{noprop}}$ is the the worst-case session i delay computed for the same characterization of internal traffic when propagation delays are zero. Also, if the number of bits in “flight” on a link (l, m) is at most

$$q_{l,m} = r_l d_{l,m}, \quad (21)$$

then

$$Q_i^* \leq \sum_{m=1}^K q_{m-1,m} + Q_i^{*,\text{noprop}}.$$

10 PGPS networks

When packet sizes are small so that maximum packet transmission time at any link of the network is negligible, we may conclude from Theorem 2 of [8], that the behavior of GPS and PGPS are (essentially) identical. Thus in this case, all of the bounds for GPS

networks in Sections 7 and 8 apply to PGPS networks as well. We now consider the more general case in which packet sizes are not negligible, and outline how our results can be extended to this case. The details of this Section are contained in [7].

The first issue of significance is that a packet is not considered to have arrived until its *last bit* has arrived. This is required for networks with heterogeneous link speeds. Thus, if $m-1$ and m are successive nodes on a session i 's route, we cannot assume, as we did in Section 8, that $S_i^{m-1} = A_i^m$. In fact, for $P(i) = \{1, 2, \dots, K_i\}$:

$$S_i^{m-1}(0, t) \geq A_i^m(0, t) \geq S_i^{m-1}(0, t) - L_i, \quad m = 2, \dots, K_i,$$

where $L_i \leq \sigma_i$ is the maximum packet size for each session i . This important difference notwithstanding, the results for $L_i = 0$, are still very useful in the more general case of non-negligible packet sizes. In particular, if we continue to assume GPS, but with non-negligible packet sizes, we can show that the internal characterization of the traffic remains unchanged, i.e., $A_i^m \sim (\sigma_i^m, \rho_i)$ for all i and $m \in P(i)$. In the [7] we use the GPS system with non-negligible packets as an intermediate step in the analysis of PGPS networks.

When the service discipline is PGPS, entire packets are served at a time, and the procedure for characterizing internal traffic of a CRST network is modified somewhat, but not in any interesting way. The next step, namely of analyzing delay along the session i route is somewhat more involved and details are contained in [7].

Theorem 5 *For each session i :*

$$D_i^{*,\text{PGPS}} \leq \tau_i^* + \sum_{m=1}^K \frac{L_{\max}}{r^m}, \quad (22)$$

where

$$\tau_i^* = \max_{\tau \geq 0} \left\{ \min\{t : G_i^K(t) = \hat{A}_i(0, \tau) + (K-1)L_{\max}\} - \tau \right\}.$$

Note that as the link speeds become faster, i.e., as $r^m \rightarrow \infty$,

$$D_i^{*,\text{PGPS}} = D_i^{*,\text{GPS}}.$$

In the remainder of this section we interpret the results of the previous section for a special CRST assignment. Under RPPS Networks $\phi_i^m = \rho_i$ for every session i and $m \in I(m)$. Recall that in Section 4 we analyzed RPPS networks when the packet sizes are negligible, and derived the bounds (10) and (11) for delay and backlog respectively. Here the corresponding bounds for PGPS service are derived.

Applying the fact that the slope of G_i^K is never less than ρ_i for each session i to (22), we have:

$$D_i^{*,\text{PGPS}} \leq \frac{\sigma_i + (K-1)L_{\max}}{\rho_i} + \sum_{m=1}^K \frac{L_{\max}}{r^m}. \quad (23)$$

The first term on the RHS is likely to dominate in most instances. In particular, in high speed networks we assume that $r^m \rightarrow \infty$, and we have

$$D_i^{*,\text{PGPS}} \leq \frac{\sigma_i + (K-1)L_{\max}}{\rho_i}. \quad (24)$$

Also, as $L_{\max} \rightarrow 0$, we get (11).

The extra delay of $\frac{(K-1)L_{\max}}{\rho_i}$ in (23) does not diminish with increasing link speed. However, as the following example shows, this term is not superfluous, but is a consequence of the PGPS service discipline: Consider a PGPS network with a large number of identically characterized sessions—i.e. $A_j \sim (\sigma, \rho)$, $\phi_j = 1$ for each session j , and all the packets have the same length, L . Every link operates at rate r , is shared by N sessions, and

$$N\rho = r - \epsilon, \quad \epsilon > 0, \quad \epsilon \approx 0. \quad (25)$$

We focus on a session i route that consists of nodes $1, 2, \dots, K$, and follow the progress of a session i packet, p , along this route. If p arrives at a node 1 at time t_1 , then assume that every other session contending for service at that node sends a packet at time t_1^- . Under PGPS, all $N-1$ packets will be served before p at node 1. Similarly, letting t_m be the time at which p arrives at node m , $2 \leq m \leq K$, we stipulate that for every other session contending for service at that node a packet arrives at time t_m^- . The delay incurred by p from these packets at node m is $\frac{(N-1)L}{r}$, which is $\approx \frac{L}{\rho}$ for large N . Thus, over all nodes in the route, this delay is $\approx \frac{KL}{\rho}$ for large N . Now letting r and L approach ∞ together, we observe that the delay term is unchanged as long as (25) continues to hold. If $L = \sigma$, the worst-case packet delay for session i will be at least $\frac{KL}{\rho}$ for large N , which corresponds to (24).

This example and (24) strongly indicate that small packet lengths should be chosen in RPPS networks so that the term $\frac{L_i}{\rho_i}$ is small. For ATM networks, in which the packets are about 400 bits long, this holds for most kinds of applications. Finally, note that the phenomenon described in our example occurs in other non-preemptive service disciplines such as FCFS as well.

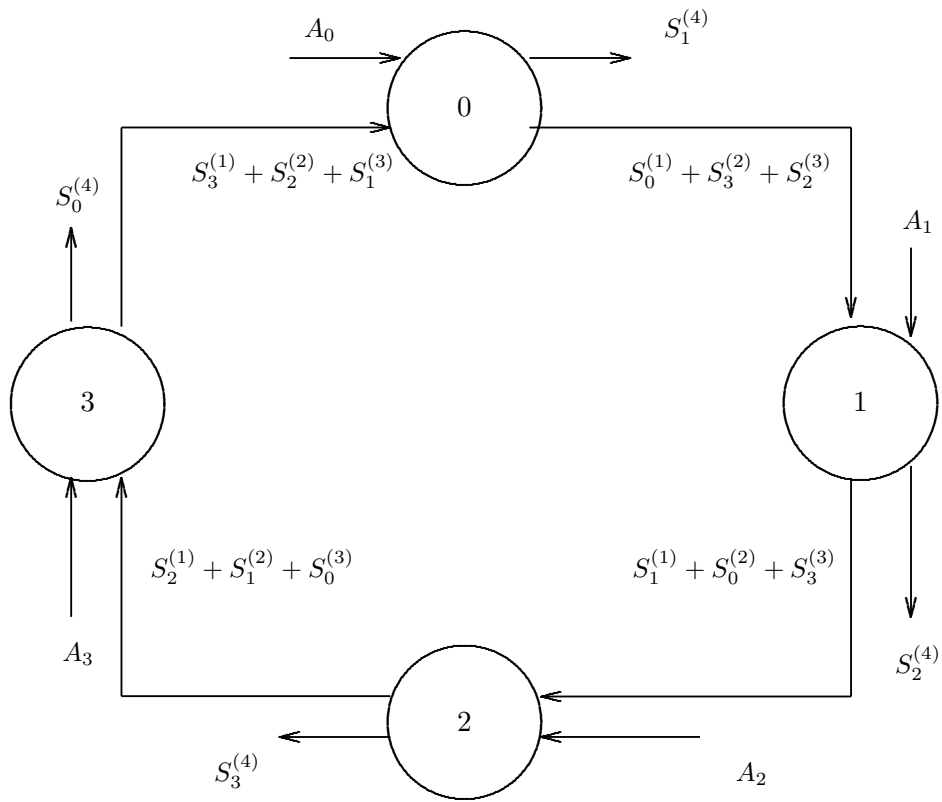
11 Conclusions

Per-session bounds were derived for the leaky bucket constrained sessions of arbitrary topology GPS and PGPS networks. With this analysis, we have provided framework for rate-based flow control in which real-time guarantees can be made to a wide variety of co-existing session types. An important part of any flow control scheme, and one that is missing from this paper is call-admission. Another area for future research is the incorporation of traffic types that require

real-time performance but that cannot predict the exact values of their leaky bucket parameters at session set-up time.

References

- [1] Rene L. Cruz. A calculus for network delay, Part I: Network elements in isolation. *IEEE Transactions on Information Theory*, 37(1):114–131, January 1991.
- [2] Rene L. Cruz. A calculus for network delay, Part II: Network analysis. *IEEE Transactions on Information Theory*, 37(1):132–141, January 1991.
- [3] Alan Demers, Srinivasan Keshav, and Scott Shenkar. Analysis and simulation of a fair queueing algorithm. *Proceedings of SIGCOMM '89*, pages 1–12, September 1989.
- [4] R. Kolla and B. Serf. The virtual feedback problem in hierarchical representations of combinatorial circuits. *Acta Informatica*, 28(5):463–476, May 1991.
- [5] C. Lu and P. R. Kumar. Distributed scheduling based on due dates and buffer prioritization. Technical report, University of Illinois Technical Report, 1990.
- [6] A. K. Parekh. *A Generalized Processor Sharing Approach to Flow Control in Integrated Services Networks*. PhD thesis, Department of Electrical Engineering and Computer Science, MIT, February 1992.
- [7] A. K. Parekh and R. G. Gallager. A Generalized Processor Sharing approach to flow control—The Multiple Node Case. Technical Report 2076, Laboratory for Information and Decision Systems, MIT, 1991.
- [8] A. K. Parekh and R. G. Gallager. A Generalized Processor Sharing approach to flow control in Integrated Services Networks—The Single Node Case. Technical Report 2040, Laboratory for Information and Decision Systems, MIT, 1991, An abridged version appears in the Proceedings of INFOCOM '92.
- [9] J. R. Perkins and P. R. Kumar. Stable distributed real-time scheduling of flexible manufacturing systems. *IEEE Transactions on Automatic Control*, AC-34(2):139–148, February 1989.



Server 0

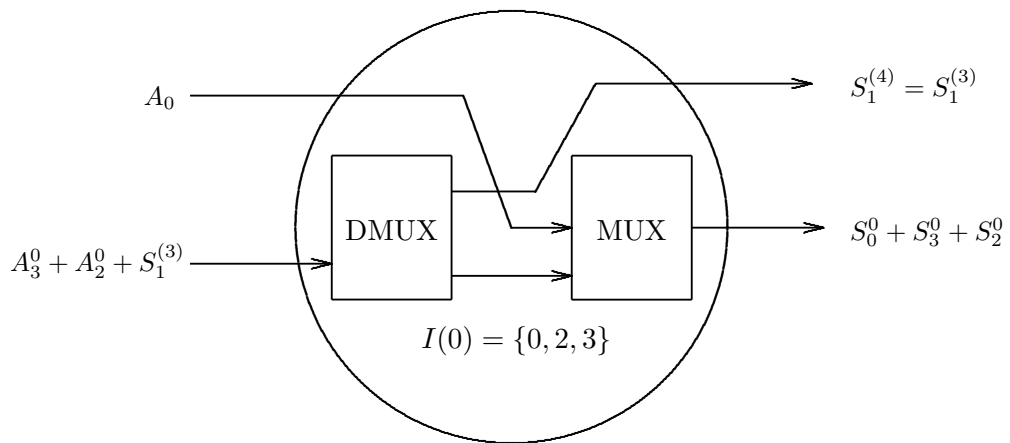
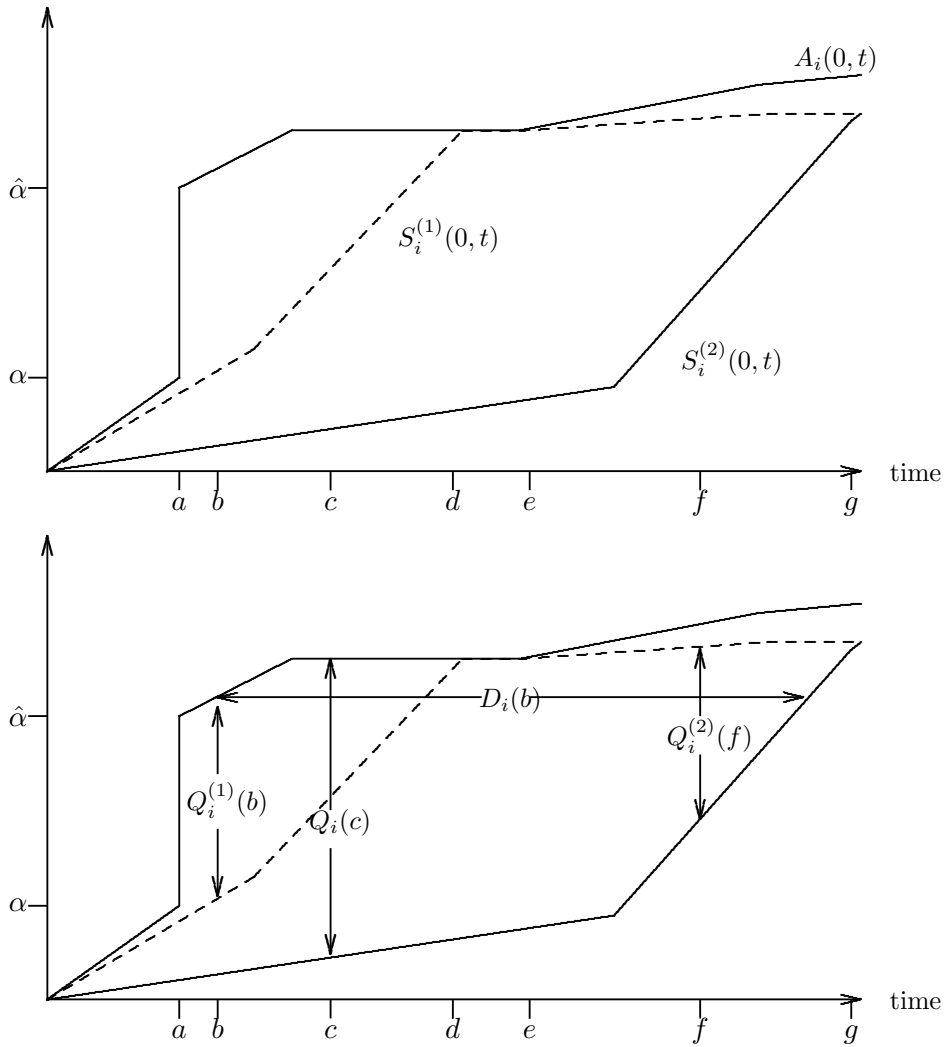
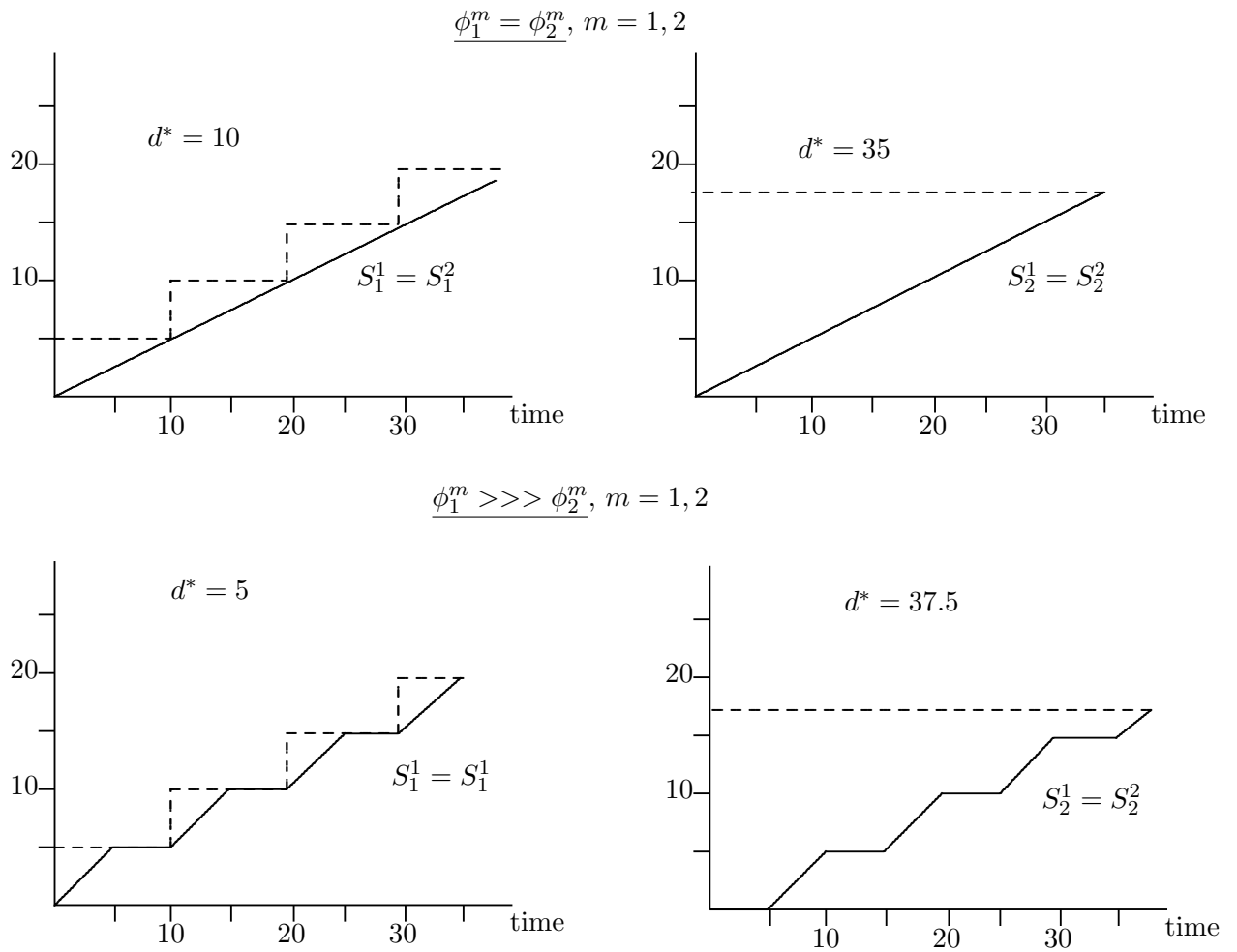


Figure 1: A four server network. The demultiplexer works instantaneously.



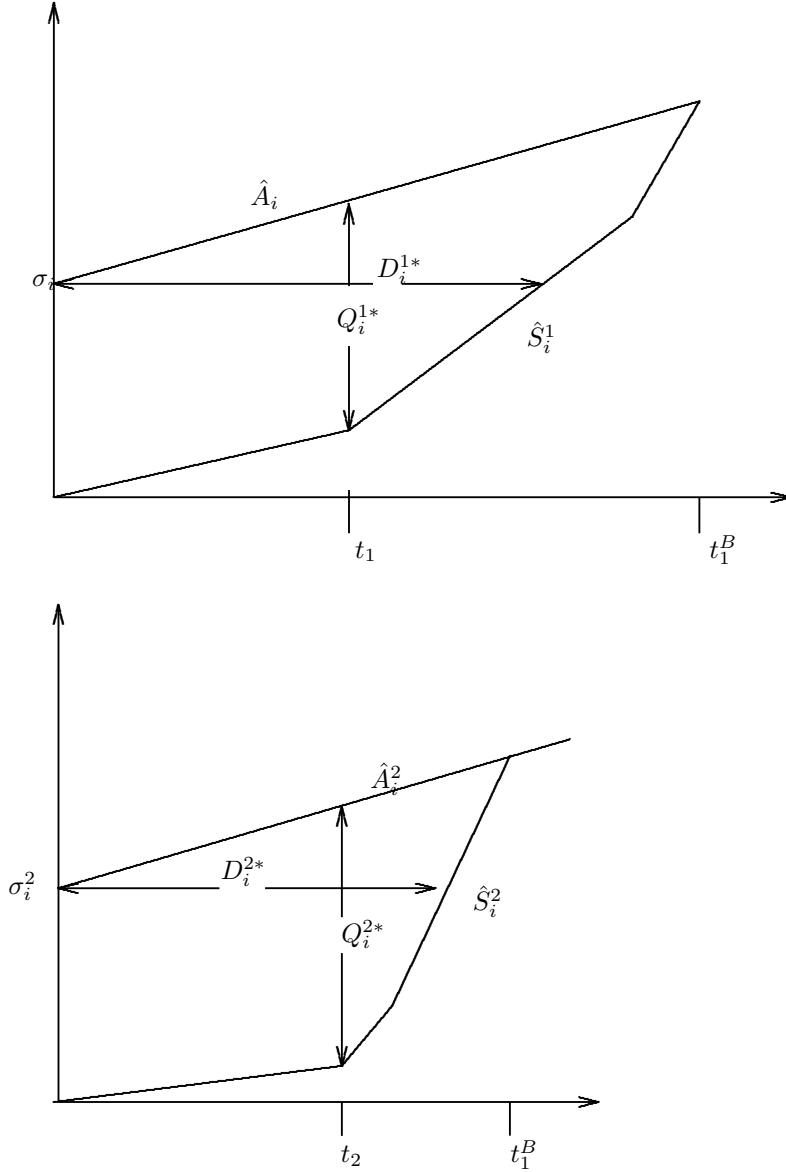
The first figure shows how session i traffic progresses through the nodes of its route. Notice that the arrival function to node 2 is the session i service function of node 1. The second figure shows how the backlog and delay can be measured and illustrates the definitions of Section 3.

Figure 2: An Example of Session i flow when $K_i = 2$.



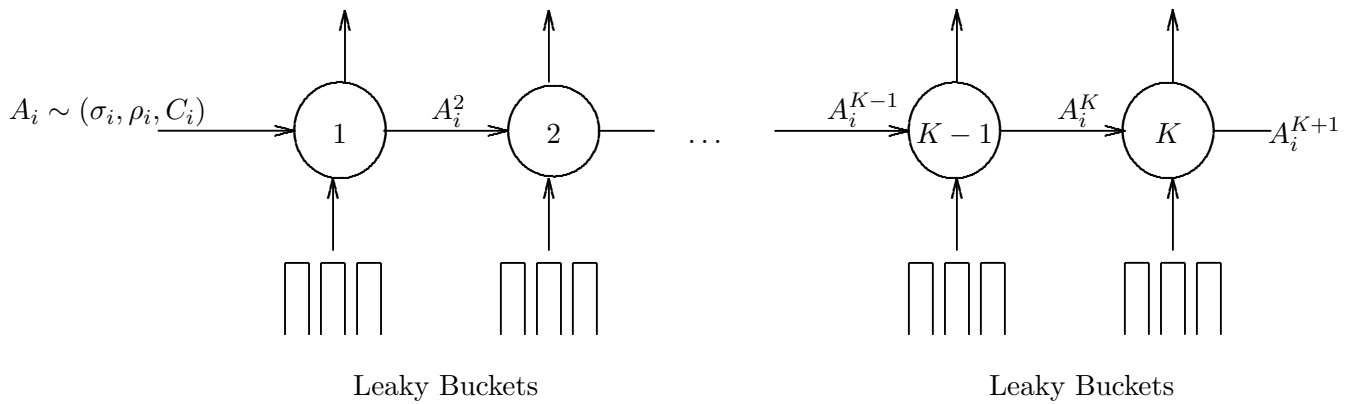
The backlog clearing rate for the Session 2 is infinitesimal when $\phi_1^m \gg \phi_2^m, m = 1, 2$. However, its delay is not increased significantly.

Figure 3: Giving delay sensitive, steady sessions large values of ϕ .



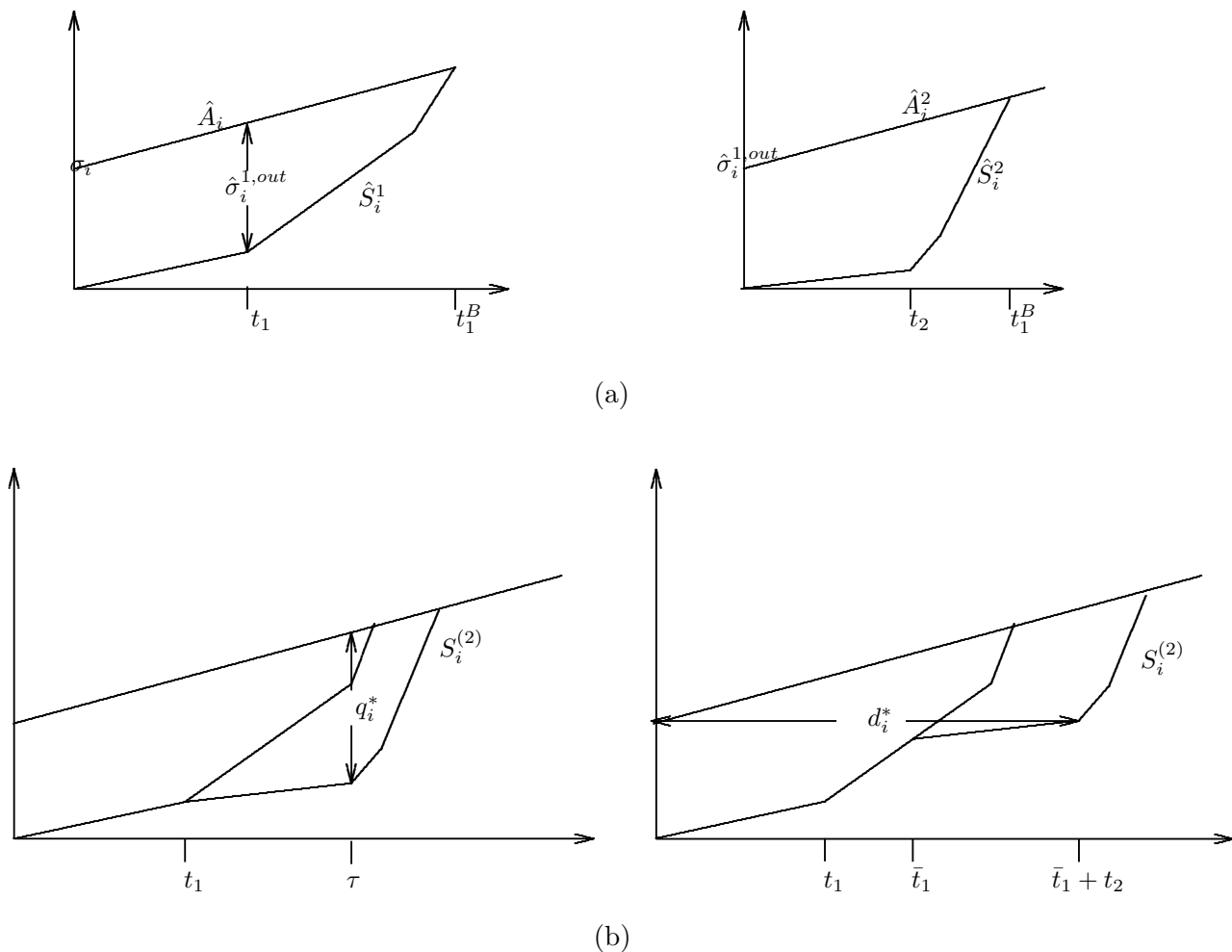
The figures (a) and (b) can be determined independently. The Additive Method yields bounds $Q_i^* \leq Q_i^{1*} + Q_i^{2*}$ and $D_i^* \leq D_i^{1*} + D_i^{2*}$.

Figure 4: The Additive Method for session i when $P(i) = \{1, 2\}$.



Session i traffic enters the network so that it is consistent with (σ_i, ρ_i) , and $A_i^k = S_i^{k-1}$ for $k = 2, 3, \dots, K$. The independent sessions at node k are free to send traffic in any manner as long as $A_j^k \sim (\sigma_j^k, \rho_j)$ for every session $j \in I(k) - \{i\}$, $k = 1, 2, \dots, K$.

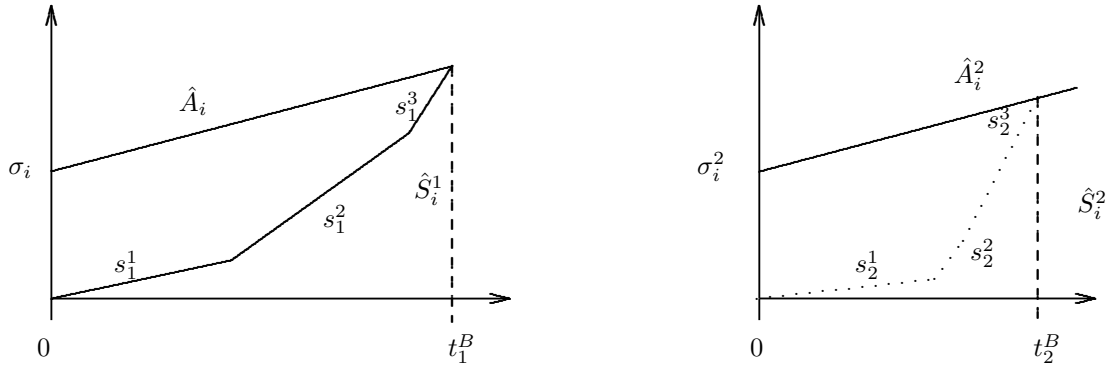
Figure 5: Analyzing the Session i route as a whole, under the Independent Sessions Relaxation.



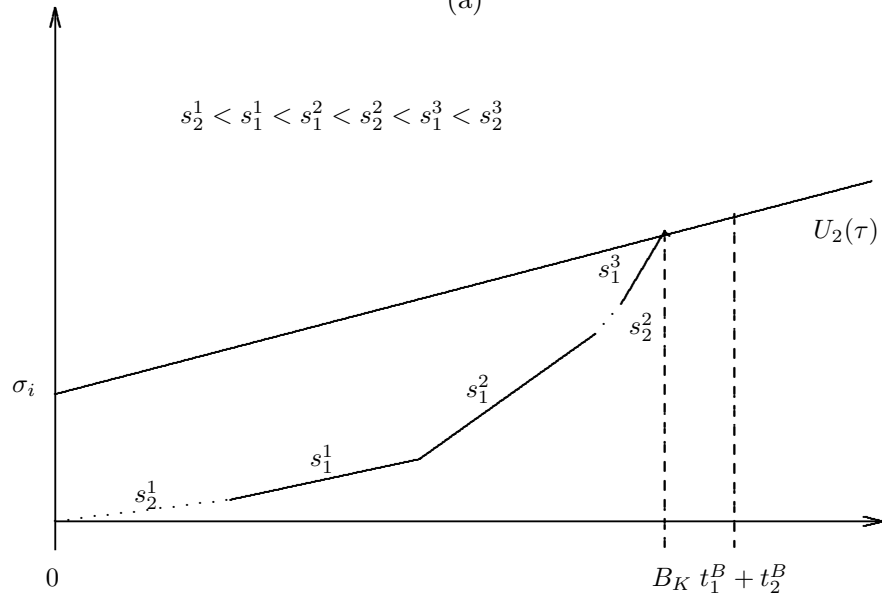
The curves \hat{S}_i^1 and \hat{S}_i^2 are shown in (a). Note that $\sigma_i^2 = \sigma_i^{1,out}$, and so \hat{S}_i^1 and \hat{S}_i^2 cannot be determined independently.

Figure (b) shows two staggered greedy regimes. In the first, the sessions in $I(2) - \{i\}$ become greedy at time t_1 , which yields a maximum backlog of q_i^* at time τ . In the second staggered greedy regime, the sessions at $I(2) = \{i\}$ wait until time \bar{t}_1 to become greedy—this results in a maximum delay of d_i^* for session i at time zero.

Figure 6: Two Staggered Greedy Regimes when $P(i) = \{1, 2\}$



(a)



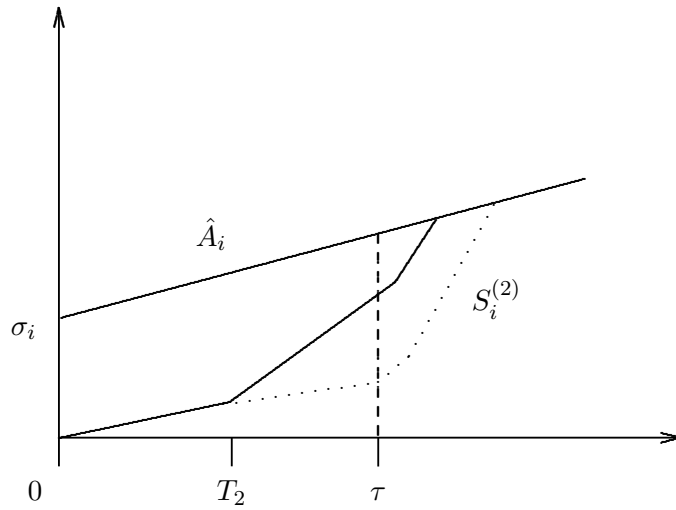
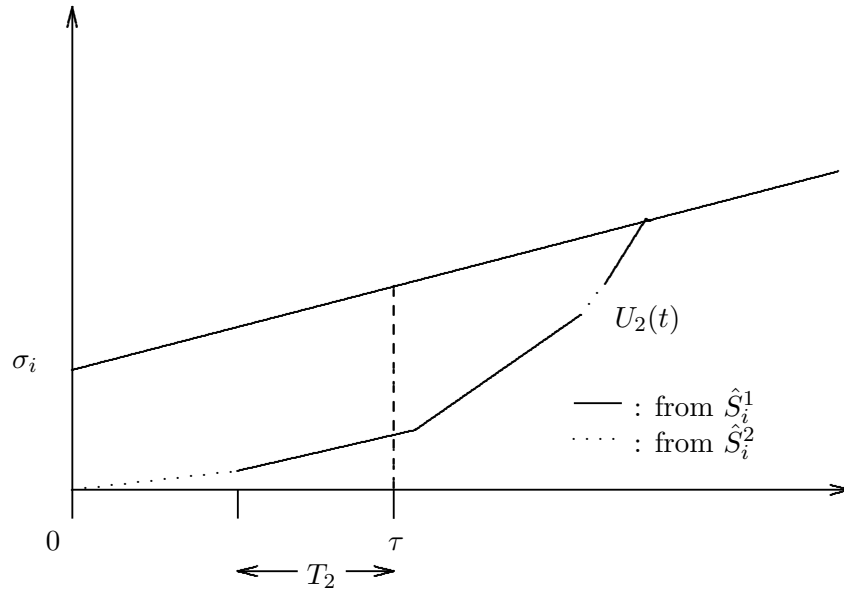
(b)

The two service curves in (a) show \hat{S}_i^1 and \hat{S}_i^2 . In (b) the line segments that make up these curves are concatenated to make a piece-wise linear convex curve that meets \hat{A}_i at time B_K . Thus

$$U_i(t) = \begin{cases} G_i^K(t) & t \leq B_K \\ \hat{A}_i(0, t) & t > B_K. \end{cases}$$

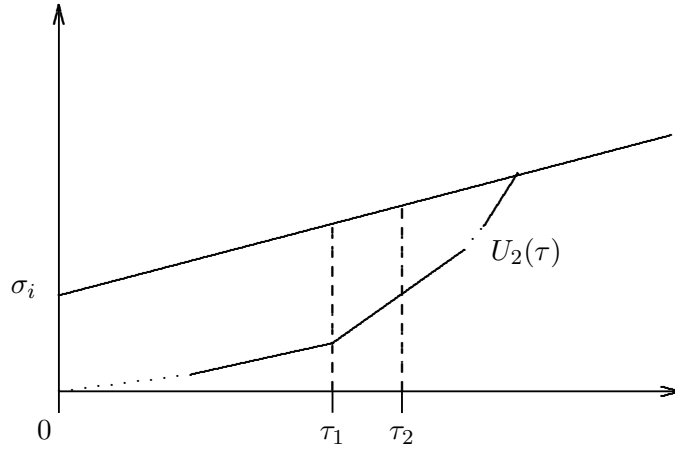
Note that the line segment with slope s_2^3 is never used in the construction of U_2 , i.e. $T_2 < t_1^B + t_2^B$.

Figure 7: An example of how U_i is constructed for $K = 2$.

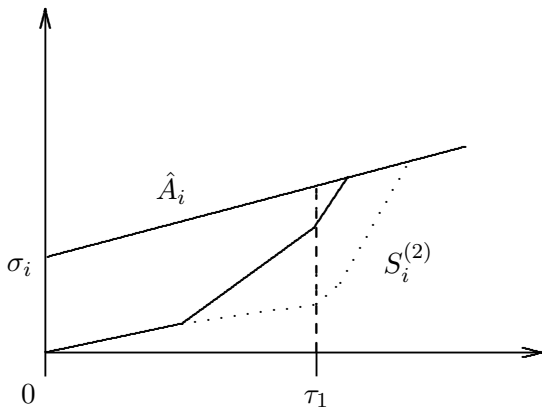


The top figure the curve U_2 that was constructed from \hat{S}_i^1 and \hat{S}_i^2 . In order to find the $(2, \tau)$ -staggered greedy regime, add the durations of the line segments taken from \hat{S}_i^1 that are in $U_2(t)$, $t \leq \tau$. This sum is T_2 , the time that the independent sessions at node 2 become greedy. This characterizes the staggered greedy regime which is shown in the bottom figure.

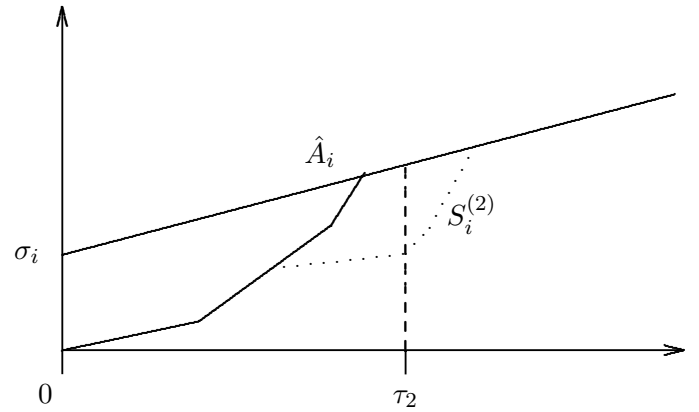
Figure 8: Computing a (k, t) -Staggered Greedy Regime when $P(i) = \{1, 2\}$



(a)



$(2, \tau_1)$ -staggered greedy



(b)

$(2, \tau_2)$ -staggered greedy

Figure (a) shows the session i universal curve. Notice that for this curve “backlog” is maximized at time τ_1 and “delay” is maximized at time τ_2 . Figure (b) shows the two staggered greedy regimes corresponding to these times. Notice that the backlog at time τ_1 in the first regime is exactly equal to the “backlog” at time τ_1 in (a), and similarly the delay at time τ_2 in the second regime is exactly equal to the “delay” at that time in (a).

Figure 9: The staggered greedy regimes that maximize backlog and delay under the independent sessions relaxation.