

A Calculus for End-to-end Statistical Service Guarantees ^{*}

Technical Report: University of Virginia, CS-2001-19

Jörg Liebeherr ^{}*

*Stephen Patek ^{**}*

Almut Burchard [†]

^{*} Department of Computer Science

^{**} Department of Systems and Information Engineering

[†] Department of Mathematics

University of Virginia

Charlottesville, VA 22904

Abstract

The deterministic network calculus offers an elegant framework for determining delays and backlog in a network with deterministic service guarantees to individual traffic flows. A drawback of the deterministic network calculus is that it only provides worst-case bounds. Here we present a network calculus for statistical service guarantees, which can exploit the statistical multiplexing gain of sources. We introduce the notion of an *effective service curve* as a probabilistic bound on the service received by an individual flow, and construct an effective service curve for a network where capacities are provisioned exclusively to aggregates of flows. Numerical examples demonstrate that the calculus is able to extract a significant amount of multiplexing gain in networks with a large number of flows.

Key Words: Quality-of-Service, Service Differentiation, Network Calculus.

1 Introduction

The deterministic network calculus recently evolved as a theory for deterministic queueing systems, and has provided powerful tools for reasoning about delay and backlog in a network with service guarantees to individual traffic flows. Using the notion of arrival envelopes and service curves [11], several recent works have shown that delay and backlog bounds can be concisely expressed in a min-plus algebra [1, 5, 8].

However, the deterministic view of traffic only provides worst-case bounds and does not take advantage of statistical multiplexing gain. The problem of trying to exploit the resource savings of statistical multiplexing while preserving the elegant formalism of the network calculus has been the subject of several studies. Kurose [16] uses the concept of stochastic ordering and obtains bounds on the distribution of delay and buffer occupancy of a flow in a network with FIFO scheduling. Chang [7] presents probabilistic bounds on output burstiness, backlog and delays in a network where the moment generating functions of arrivals are exponentially bounded. Different bounds for exponentially bounded arrivals are derived by Yaron and

^{*}This work is supported in part by the National Science Foundation through grants ANI-9730103, ECS-9875688 (CAREER), ANI-9903001, DMS-9971493, and ANI-0085955, and by an Alfred P. Sloan research fellowship.

Sidi [23] and Starobinski and Sidi [22]. Results on statistical end-to-end delay guarantees in a network have been obtained for specific scheduling algorithms, such as EDF [20, 21], and GPS [13], and a class of co-ordinated scheduling algorithms [2, 17]. Several researchers have considered probabilistic formulations of service curves. Cruz defines a probabilistic service curve which violates a given deterministic service curve according to a certain distribution [12]. Chang (see [9], Chp. 7) presents a statistical network calculus for ‘dynamic F-servers’. Finally, Knightly and Qiu [19] derive ‘statistical service envelopes’ as time-invariant lower bounds on the service received by an aggregate of flows.

This paper proposes a network calculus for statistically multiplexed traffic, expressed in the min-plus algebra, where network capacities are allocated to aggregates of flows. This is different from the per-flow capacity allocation generally applied in the deterministic network calculus. Within this context, we define an *effective service curve*, which is, with high certainty, a bound on the service received by a single flow. So, we will consider probabilistic per-flow service guarantees for networks where resources are reserved for aggregates. We will show that the main results of the deterministic network calculus carry over to the statistical framework we present.

The results in this paper are set in a continuous time model with fluid left-continuous traffic arrival functions, as is common for network delay analysis in the deterministic network calculus. We refer to [9] for the issues involved in relaxing these assumptions for the analysis of packet networks. A node represents a router (or switch) in a network. The transmission rate at a node corresponds to the capacity of an output link of a router. Packetization delays and other effects of discrete-sized packets, such as the non-preemption of packet transmission, are ignored. When analyzing delays in a network, all processing overhead and propagation delays are ignored.

In the numerical examples presented in this paper, we assume a ‘regulated adversarial traffic’ model where (1) arrivals from each flow into the network are constrained by a deterministic regulator and (2) traffic arrivals from different flows are statistically independent. The regulated adversarial traffic model has been used by several researchers, e.g., [14, 15], for modeling aggregates of sources, which are policed or shaped, but for which arrival distributions are not readily available.

The remaining sections of this paper are structured as follows. In Section 2, we review the notation and key results of the deterministic network calculus. In Section 3 we introduce effective service curves and present the results for a statistical network calculus in terms of effective service curves. In Section 4 we show how to construct effective service curves for individual flows at a node where service is allocated to an aggregate of flows. In Section 5, we show how to build ‘effective envelopes’ [4], which are used in our construction of effective service curves. In Section 6, we discuss numerical examples for single node and multi-node networks and evaluate the statistical multiplexing gain achievable with effective service curves.

2 Network Calculus Preliminaries

The deterministic network calculus provides concise expressions for upper bounds on the backlog and delay experienced by an individual flow at one or more network nodes. An attractive feature of the network calculus is that end-to-end bounds can often be easily obtained from manipulations of the per-node bounds.

In this section we review some notation and results from the deterministic network calculus, as needed later in the paper. However, this section is not a comprehensive summary of the network calculus. For a complete discussion we refer to [1, 6, 9].

2.1 Operators

Much of the formal framework of the network calculus can be elegantly expressed in a min-plus algebra [3], complete with convolution and deconvolution operators for functions. Generally, the functions in this paper are non-negative and left-continuous, defined over time intervals $[0, t]$. We assume for a given function f that $f(t) = 0$ if $t < 0$.

The *convolution* $f * g$ of two functions f and g , is defined as

$$f * g(t) = \inf_{\tau \in [0, t]} \{f(t - \tau) + g(\tau)\} . \quad (1)$$

The *deconvolution* $f \oslash g$ of two functions f and g , is defined as

$$f \oslash g(t) = \sup_{\tau \geq 0} \{f(t + \tau) - g(\tau)\} . \quad (2)$$

We refer to [3, 6, 9] for a detailed discussion of the properties of the min-plus algebra and the properties of the convolution and deconvolution operators.

2.2 Arrival functions and Service Curves

Let us consider the traffic arrivals to a single network node. The arrivals of a flow in the time interval $[0, t]$ are given in terms of a function $A(t)$. The departures of a flow from the node in the time interval $[0, t]$ are denoted by $D(t)$, with $D(t) \leq A(t)$. The backlog of a flow at time t , denoted by $B(t)$, is given by

$$B(t) = A(t) - D(t) . \quad (3)$$

The delay at time t , denoted as $W(t)$, is the delay experienced by an arrival which departs at time t , given by

$$W(t) = \inf\{d \geq 0 \mid A(t - d) \leq D(t)\} . \quad (4)$$

If arrival and departure functions are plotted as functions of time, then $B(t)$ and $W(t)$, respectively, are the vertical and horizontal differences between arrival and departure functions. We will use $A(x, y)$ and $D(x, y)$ to denote the arrivals and departures in the time interval $[x, y]$, with $A(x, y) = A(y) - A(x)$ and $D(x, y) = D(y) - D(x)$.

We have the following assumptions on the arrival functions.

(A1) *Non-Negativity*. The arrivals in any interval of time are non-negative. That is, for any $x < y$, we have $A(y) - A(x) \geq 0$.

(A2) *Upper Bound*. The arrivals A of a flow are bounded by a deterministic subadditive function A^* , called the *arrival envelope*,¹ such that $A(t + \tau) - A(t) \leq A^*(\tau)$ for all $t, \tau \geq 0$.²

¹ A function E is called an *envelope* for a function f if $f(t + \tau) - f(t) \leq E(\tau)$ for all $t, \tau \geq 0$, or, equivalently, if $f(t) \leq E * f(t)$, for all $t \geq 0$.

² A function f is subadditive if $f(x + y) \leq f(x) + f(y)$, for all $x, y \geq 0$, or, equivalently, if $f(t) = f * f(t)$.

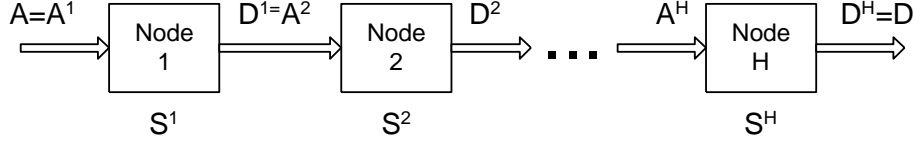


Figure 1: Traffic of a flow through a set of H nodes. Let A^h and D^h denote the arrival and departures at the h -th node, with $A^1 = A$, $A^h = D^{h-1}$ for $h = 2, \dots, H$ and $D^H = D$.

The upper bound given by A^* assumes that the traffic of a flow is policed or shaped by a traffic conditioning function, such as a leaky bucket.

The service guaranteed to a flow is expressed in terms of service curves. A *minimum service curve* for a flow is a function S which specifies a lower bound on the service given to a flow such that, for all $t \geq 0$,

$$D(t) \geq A * S(t) . \quad (5)$$

A *maximum service curve* for a flow is a function \bar{S} which specifies an upper bound on the service given to a flow such that, for all $t \geq 0$,

$$D(t) \leq A * \bar{S}(t) . \quad (6)$$

The following theorem summarizes some key results of the deterministic network calculus. These results have been derived in [1, 5, 8]. We follow the notation used in [1]. A proof of the theorem is included in Appendix A.

Theorem 1 Deterministic Network Calculus. *Given a flow with arrival envelope A^* and with minimum and maximum service curves S and \bar{S} , the following hold:*

1. **Output Envelope:** *The function $D^* = A^* \oslash S$ is an envelope for the departures, in the sense that, for all $t, \tau \geq 0$,*

$$D^*(t) \geq D(t + \tau) - D(\tau) . \quad (7)$$

2. **Backlog Bound:** *An upper bound for the backlog, denoted by b_{max} , is given by*

$$b_{max} = A^* \oslash S(0) . \quad (8)$$

3. **Delay Bound:** *An upper bound for the delay, denoted by d_{max} , is given by*

$$d_{max} = \inf \{d \geq 0 \mid \forall t \geq 0 : A^*(t - d) \leq S(t)\} . \quad (9)$$

4. **Network Service Curve:** *Suppose a flow passes through H nodes in series, as shown in Figure 1, and suppose the flow is offered minimum and maximum service curves S^h and \bar{S}^h , respectively, at each node $h = 1, \dots, H$. Then, the sequence of nodes provides minimum and maximum service curves S^{net} and \bar{S}^{net} , which are given by*

$$S^{net} = S^1 * S^2 * \dots * S^H , \quad (10)$$

$$\bar{S}^{net} = \bar{S}^1 * \bar{S}^2 * \dots * \bar{S}^H . \quad (11)$$

S^{net} and \bar{S}^{net} will be referred to as network service curves. With Theorem 1, network service curves can be used to determine bounds on delay and backlog in a network. There are many additional properties and refinements that have been derived for the deterministic calculus. However, in this paper we will concern ourselves only with the results above.

3 Statistical Network Calculus

A limitation of the deterministic network calculus is that the deterministic view of traffic only yields pessimistic worst-case bounds, which do not take advantage of the statistical multiplexing of flows when multiple flows are carried over the same link. We will now approach the network calculus in a probabilistic framework. Arrivals and departures from a flow to the network in the time interval $[0, t)$ are described by random processes $A(t)$ and $D(t)$. The random processes are defined over some probability space that we suppress in our notation. The statistical network calculus makes service guarantees for individual flows, where each flow is allocated a probabilistic service in the form of an ‘effective service curve’.

Definition 1 *Given a flow with arrival process A , which satisfies assumptions (A1)–(A2), a (minimum) effective service curve is a function \mathcal{S}^ε that satisfies for all $t \geq 0$,*

$$Pr\{D(t) \geq A * \mathcal{S}^\varepsilon(t)\} \geq 1 - \varepsilon, \quad (12)$$

and a maximum effective service curve is a function $\bar{\mathcal{S}}^\varepsilon$ that satisfies

$$Pr\{D(t) \leq A * \bar{\mathcal{S}}^\varepsilon(t)\} \geq 1 - \varepsilon. \quad (13)$$

Note that in this definition, an effective service curve is a non-random function. The following theorem states some key results for the network calculus in terms of effective service curves.

Theorem 2 Statistical Network Calculus for Flows. *Given the arrival process A of a flow with arrival envelope A^* and given effective service curves \mathcal{S}^ε and $\bar{\mathcal{S}}^\varepsilon$, the following hold:*

1. **Output Envelope:** *The function $A^* \odot \mathcal{S}^\varepsilon$ is a probabilistic bound for the departures, in the sense that, for all $t, \tau \geq 0$,*

$$Pr\{D(t, t + \tau) \leq A^* \odot \mathcal{S}^\varepsilon(\tau)\} \geq 1 - \varepsilon. \quad (14)$$

2. **Backlog Bound:** *A probabilistic bound for the backlog is given by $b_{max} = A^* \odot \mathcal{S}^\varepsilon(0)$, in the sense that, for all $t \geq 0$,*

$$Pr\{B(t) \leq b_{max}\} \geq 1 - \varepsilon. \quad (15)$$

3. **Delay Bound:** *A probabilistic bound for the delay is given by,*

$$d_{max} = \inf\{d \geq 0 \mid \forall t \geq 0 : A^*(t - d) \leq \mathcal{S}^\varepsilon(t)\}, \quad (16)$$

in the sense that, for all $t \geq 0$,

$$Pr\{W(t) \leq d_{max}\} \geq 1 - \varepsilon. \quad (17)$$

4. **Network Service Curve:** Suppose the flow passes through H network nodes in series and is offered minimum and maximum service curves $\mathcal{S}^{h, \varepsilon_h}$ ($\overline{\mathcal{S}}^{h, \varepsilon_h}$), respectively, at each node $h = 1, \dots, H$. Then, minimum and maximum effective network service curves are given by

$$\mathcal{S}^{net, \varepsilon_1 + \varepsilon_2 + \dots + \varepsilon_H} = \mathcal{S}^{1, \varepsilon_1} * \mathcal{S}^{2, \varepsilon_2} * \dots * \mathcal{S}^{H, \varepsilon_H}, \quad (18)$$

$$\overline{\mathcal{S}}^{net, \varepsilon_1 + \varepsilon_2 + \dots + \varepsilon_H} = \overline{\mathcal{S}}^{1, \varepsilon_1} * \overline{\mathcal{S}}^{2, \varepsilon_2} * \dots * \overline{\mathcal{S}}^{H, \varepsilon_H}. \quad (19)$$

That is, for all $t \geq 0$,

$$Pr \{D(t) \geq A^1 * (\mathcal{S}^{1, \varepsilon_1} * \mathcal{S}^{2, \varepsilon_2} * \dots * \mathcal{S}^{H, \varepsilon_H})(t)\} \geq 1 - (\varepsilon_1 + \varepsilon_2 + \dots + \varepsilon_H), \quad (20)$$

$$Pr \{D(t) \leq A^1 * (\overline{\mathcal{S}}^{1, \varepsilon_1} * \overline{\mathcal{S}}^{2, \varepsilon_2} * \dots * \overline{\mathcal{S}}^{H, \varepsilon_H})(t)\} \geq 1 - (\varepsilon_1 + \varepsilon_2 + \dots + \varepsilon_H). \quad (21)$$

The proof of the theorem is given in Appendix B. In the next section we show how effective service curves can make probabilistic statements about service guarantees for individual flows in a network where deterministic service curves \mathcal{S}_C (and $\overline{\mathcal{S}}_C$) are provisioned to flow aggregates.

4 Construction of Effective Service Curves

In this section, we present the construction of a minimum effective service curve for a single flow in a network where bandwidth is allocated to aggregates of flows. The effective service curve for a given flow is determined from the unused bandwidth that is allocated to the aggregate of flows. The construction of the effective service curve is based on the notion of effective envelopes from [4].

4.1 Effective Envelopes

Let \mathcal{C} denote a set of flows. The arrival and departure processes for each flow $j \in \mathcal{C}$ will be denoted by A_j and D_j , respectively. We use $A_{\mathcal{C}}$ and $D_{\mathcal{C}}$ to denote the aggregate arrivals and departures from class \mathcal{C} at a network node, that is,

$$A_{\mathcal{C}}(t) = \sum_{j \in \mathcal{C}} A_j(t), \text{ and} \quad (22)$$

$$D_{\mathcal{C}}(t) = \sum_{j \in \mathcal{C}} D_j(t). \quad (23)$$

A deterministic arrival envelope for the aggregate is

$$A_{\mathcal{C}}^*(t) = \sum_{j \in \mathcal{C}} A_j^*(t). \quad (24)$$

The backlog and delay, respectively, for the set of flows are defined by

$$B_{\mathcal{C}}(t) = A_{\mathcal{C}}(t) - D_{\mathcal{C}}(t), \text{ and} \quad (25)$$

$$W_{\mathcal{C}}(t) = \inf\{d \geq 0 \mid A_{\mathcal{C}}(t - d) \leq D_{\mathcal{C}}(t)\}. \quad (26)$$

With the above notation, we now define a number of probabilistic bounds on the arrivals for an aggregate set of flows. These bounds are called *effective envelopes*. In Definition 2, we recall the definitions of effective envelopes from [4].

Definition 2 Given a set \mathcal{C} of flows that satisfy assumptions (A1)–(A2).

1. A local effective envelope for $A_{\mathcal{C}}$ is a function $\mathcal{G}_{\mathcal{C}}^{\varepsilon}$ such that for all t and τ

$$\Pr\left\{A_{\mathcal{C}}(t, t + \tau) \leq \mathcal{G}_{\mathcal{C}}^{\varepsilon}(\tau)\right\} \geq 1 - \varepsilon. \quad (27)$$

2. A global effective envelope for $A_{\mathcal{C}}$ for intervals of length l is a subadditive function $\mathcal{H}_{\mathcal{C}}^{l, \varepsilon}$ such that for each interval I_l of length l ,

$$\Pr\left\{A_{\mathcal{C}}(t, t + \tau) \leq \mathcal{H}_{\mathcal{C}}^{l, \varepsilon}(\tau), \forall t, \tau : [t, t + \tau] \subseteq I_l\right\} \geq 1 - \varepsilon. \quad (28)$$

Thus, a local effective envelope provides a bound for arrivals in the time interval $[0, t]$, which is violated with probability at most ε . On the other hand, a global effective envelope is a probabilistic bound for all subintervals in any interval I_l of length l . The global effective envelope is instrumental for the construction of our proposed effective service curve. The local effective envelope is relevant as it will be used to generate a global effective envelope (see Section 5).

Remark:

- Note that the probabilistic bound for the output traffic $A^* \oslash \mathcal{S}^{\varepsilon}$ from Theorem 2 is a local effective envelope.
- The global effective envelope in the above definition is equivalent to that in [4], where a global effective envelope is defined by the property that for each interval I_l of length l

$$\Pr\left\{\mathcal{E}_{\mathcal{C}}(\tau) \leq \mathcal{H}_{\mathcal{C}}^{l, \varepsilon}(\tau), \forall \tau \leq l\right\} \geq 1 - \varepsilon, \quad (29)$$

where

$$\mathcal{E}_{\mathcal{C}}(\tau) = \sup_{[t, t + \tau] \subseteq I_l} A_{\mathcal{C}}(t, t + \tau). \quad (30)$$

4.2 Effective Service Curves

We consider a set \mathcal{C} of flows which satisfy assumptions (A1)–(A2), and construct a minimum effective service curve for flow j with arrival process A_j . We assume that the aggregate set of flows is allocated a (deterministic) minimum service curve, denoted by $S_{\mathcal{C}}$, and the set $\mathcal{C} - \{j\}$ is allocated a maximum service curve of $\overline{S}_{\mathcal{C} - \{j\}}$. Let $\mathcal{H}_{\mathcal{C} - \{j\}}^{t, \varepsilon}$ denote a global effective envelope for the arrivals from $\mathcal{C} - \{j\}$ for time intervals of length $t \geq 0$. With this notation, a minimum effective service curve for flow i is given by the next theorem.

Theorem 3 *The function*

$$\mathcal{S}_j^{\varepsilon}(t) = [S_{\mathcal{C}} - \mathcal{H}_{\mathcal{C} - \{j\}}^{t, \varepsilon} * \overline{S}_{\mathcal{C} - \{j\}}]_+(t) \quad (31)$$

is a (minimum) effective service curve for flow $j \in \mathcal{C}$.³

³We use “ $[f]_+(t) = \max\{f(t), 0\}$ ”.

The above minimum effective service curve does not assume knowledge of the scheduling algorithm used to determine the order of transmission of the aggregate of flows. Thus, the effective service curve is expected to be pessimistic for most scheduling algorithms, including FIFO. The minimum effective service curve is least conservative if flows in the set $\mathcal{C} - \{j\}$ are transmitted with higher priority than flow j .

Remarks:

- Note that $\mathcal{S}_j^\varepsilon$ as defined in Eqn. (31) need not be increasing. $\mathcal{S}_j^\varepsilon$ will be an increasing function if $\mathcal{H}_{\mathcal{C}-\{j\}}^{t,\varepsilon}$ is concave.
- The minimum effective service curve in the theorem has a corresponding version in the deterministic network calculus, which is given by $S_j(t) = [S_{\mathcal{C}} - A_{\mathcal{C}-\{j\}}^* * \overline{S}_{\mathcal{C}-\{j\}}]_+(t)$ with $A_{\mathcal{C}-\{j\}}^* = A_{\mathcal{C}}^* - A_j^*$. However, this deterministic service curve will be positive only for large values of t [6].

The following corollary states simpler bounds on the minimum effective service to flow j .

Corollary 1 *Using the same notation as in Theorem 3, and assuming that $\overline{S}_{\mathcal{C}}(t) \geq \overline{S}_{\mathcal{C}-\{j\}}(t)$, the following are (minimum) effective service curves for flow $j \in \mathcal{C}$:*

1. $\mathcal{S}_j^\varepsilon(t) = [S_{\mathcal{C}} - \mathcal{H}_{\mathcal{C}-\{j\}}^{t,\varepsilon} * \overline{S}_{\mathcal{C}}]_+(t)$,
2. $\mathcal{S}_j^\varepsilon(t) = [S_{\mathcal{C}} - \mathcal{H}_{\mathcal{C}}^{t,\varepsilon} * \overline{S}_{\mathcal{C}}]_+(t)$,
3. $\mathcal{S}_j^\varepsilon(t) = [S_{\mathcal{C}} - \mathcal{H}_{\mathcal{C}-\{j\}}^{t,\varepsilon}]_+(t)$,
4. $\mathcal{S}_j^\varepsilon(t) = [S_{\mathcal{C}} - \mathcal{H}_{\mathcal{C}}^{t,\varepsilon}]_+(t)$.

Proof. The first two service curves follow from $\overline{S}_{\mathcal{C}}(t) \geq \overline{S}_{\mathcal{C}-\{j\}}(t)$ and since $\mathcal{H}_{\mathcal{C}}^{t,\varepsilon}$ is greater than $\mathcal{H}_{\mathcal{C}-\{j\}}^{t,\varepsilon}$. The last two service curves in addition exploit that $f(t) \geq f * g(t)$, which follows from the definition of the convolution operator. \square

Thus, minimum effective service curves for single flows can be determined even if only information about the aggregate reservations to a set of flows is available. In our numerical examples, we will generally work with the last and most pessimistic effective service curve. We will show that even with these very loose bounds, we are able to extract a significant amount of the multiplexing gain if the number of flows is large.

4.3 Busy Periods for Fixed-Rate Service Curves

Instead of constructing an effective envelope $\mathcal{H}_{\mathcal{C}-\{j\}}^{t,\varepsilon}$ for each value of $t \geq 0$, as given in Theorem 2, it is more practical to select a value T large enough to cover the longest interval of interest, i.e., the longest busy period, and use $\mathcal{H}_{\mathcal{C}-\{j\}}^{T,\varepsilon}$ for all values of t . One can verify from the proof of Theorem 3 that the resulting service curve satisfies the stronger property

$$Pr(D_j(t) \geq A_j * \mathcal{S}_j^\varepsilon(t), \forall t : t \leq T) \geq 1 - \varepsilon. \quad (32)$$

We now address the issue of finding good bounds for the length of a busy period. A busy period for a set of flows \mathcal{C} is a time interval during which the backlog from the flows in \mathcal{C} remains positive. More precisely, for

any $t \geq 0$, let \underline{t}_C denote the start time and let \bar{t}_C denote the end time of the busy period of the set C around time t . We have

$$\underline{t}_C = \sup\{\tau \leq t \mid B_C(\tau) = 0\} , \text{ and} \quad (33)$$

$$\bar{t}_C = \inf\{\tau \geq t \mid B_C(\tau) = 0\} . \quad (34)$$

If the minimum service curve S_C is a fixed-rate function, that is, $S_C(t) = \text{const.} \cdot t$, a deterministic bound for the length of the longest busy period can be obtained from the deterministic calculus in Section 2 by

$$T_o = \inf\{\tau > 0 \mid A_C^*(\tau) \leq S_C(\tau)\} , \quad (35)$$

in the sense that, for all $t \geq 0$,

$$T_o \geq \bar{t}_C - \underline{t}_C . \quad (36)$$

This bound, however, is generally very conservative. The following theorem can be used to find a less conservative estimate for the length of the longest busy period.

Theorem 4 *Let T_o be given by Eqn. (35), and define recursively*

$$T_n = \inf\left\{\tau \leq T_{n-1} \mid \mathcal{H}_C^{2T_{n-1}, \varepsilon}(\tau) \leq S_C(\tau)\right\} \quad (37)$$

where S_C is a constant-rate function and where $\mathcal{H}_C^{2T_{n-1}, \varepsilon}$ is a global effective envelope for the arrivals of a set C . Then each T_n is a probabilistic bound on the busy periods of the set C , in the sense that, for any $t \geq 0$,

$$Pr\{\bar{t}_C - \underline{t}_C \leq T_n\} \geq 1 - n\varepsilon . \quad (38)$$

The theorem is proven in Appendix D. Since the bound T_n for the busy period typically decreases with n while the error $n\varepsilon$ increases, one needs to pick a ‘good’ value for n . Our examples in Section 6 show the bound T_1 for the busy period to be significantly below the deterministic bound T_o , whereas successive T_n for $n > 1$ do not result in noticeable improvements.

5 Effective Envelopes for Heterogeneous Traffic

The presentation of the effective envelopes in Section 4 does not depend on a specific arrival model, but also does not offer any guidance for constructing the envelopes. Here, we use the constructions for $\mathcal{G}_C^\varepsilon$ and $\mathcal{H}_C^{l, \varepsilon}$ from [4], adopting an adversarial traffic model [14], where arrivals of flows to the network can individually exhibit a worst-case arrival pattern as allowed by (A2), but sources do not conspire to construct a joint worst-case. In addition to assumptions (A1) and (A2) from Section 2, we assume that the following hold for the arrival processes.

(A3) *Stationarity.* The arrival processes are *stationary*, i.e., $\forall t_1, t_2 \geq 0$ we have $Pr\{A(t_1, t_1 + \tau) \leq x\} = Pr\{A(t_2, t_2 + \tau) \leq x\}$.

(A4) *Independence.* The arrivals from two flows $i, j \in C$, A_i and A_j , are stochastically independent.

We emphasize that these assumptions only hold for the arrivals to the first node of a flow’s route through the network. Since buffering and scheduling distort traffic and introduce correlations between flows, assumptions (A2)–(A4) may not hold after traffic has passed through a node.

5.1 Local Effective Envelopes for Heterogeneous Traffic

We present a construction of local effective envelopes for flows which satisfy assumptions (A1)–(A4). We allow heterogeneous flows, that is, we allow that flows can have different arrival envelopes. The following construction adapts the derivations in [4] to heterogeneous flows. The construction of effective envelopes $\mathcal{G}_{\mathcal{C}}^{\varepsilon}$ for a set \mathcal{C} of flows uses the moment generating function of A_j , denoted as $M_j(s, t) = E[e^{A_j(\tau, \tau+t)s}]$. As shown in [4], if assumptions (A1)–(A3) hold, we obtain $M_j(s, t) \leq \overline{M}_j(s, t)$, where

$$\overline{M}_j(s, t) = 1 + \frac{\rho_j t}{A_j^*(t)} (e^{s A_j^*(t)} - 1), \quad (39)$$

and where $\rho_j := \lim_{t \rightarrow \infty} A_j^*(t)/t$ is assumed to exist.

With assumption (A4) and with the bound in Eqn. (39), we obtain from the Chernoff bound that

$$Pr\{A_{\mathcal{C}}(t) \geq x\} \leq e^{-xs} \prod_{j \in \mathcal{C}} \overline{M}_j(s, t). \quad (40)$$

Setting the right hand side equal to ε and solving for x gives

$$x = \frac{1}{s} \left(\sum_{j \in \mathcal{C}} \log \overline{M}_j(s, t) + \log \varepsilon^{-1} \right). \quad (41)$$

Any choice of $s > 0$ yields a point of an effective envelope for the arrivals from \mathcal{C} . We select the value of the effective envelope at t to be

$$\mathcal{G}_{\mathcal{C}}^{\varepsilon}(t) = \inf_{s > 0} \frac{1}{s} \left(\sum_{j \in \mathcal{C}} \log \overline{M}_j(s, t) + \log \varepsilon^{-1} \right). \quad (42)$$

With this choice, $\mathcal{G}_{\mathcal{C}}^{\varepsilon}(t) \leq A_{\mathcal{C}}^*(t)$ is always satisfied. Since the derivative of the right hand side of Eqn. (41) is increasing in s , there is at most one minimum. This minimum can be found by searching for the zero of the derivative.

5.2 Global Effective Envelopes for Heterogeneous Traffic

We will construct a global effective envelope $\mathcal{H}_{\mathcal{C}}^{l, \varepsilon}$ from a local effective envelope $\mathcal{G}_{\mathcal{C}}^{\varepsilon}$, for a set \mathcal{C} of heterogeneous flows. The local effective envelope $\mathcal{G}_{\mathcal{C}}^{\varepsilon}$ may be constructed as in the previous subsection. Alternatively, $\mathcal{G}_{\mathcal{C}}^{\varepsilon}$ may be obtained from Theorem 2 as a probabilistic output envelope. If assumption (A2) holds for each flow A_j in \mathcal{C} , we have $A_{\mathcal{C}}^*(t) = \sum_{j \in \mathcal{C}} A_j^*(t)$. If assumption (A2) cannot be made, e.g., if the local effective envelope is constructed by Theorem 2, we use $A_{\mathcal{C}}^*(t) = Ct$, where C is the capacity of the link. If no deterministic bound is available, we allow $A_{\mathcal{C}}^*(t) = \infty$ for $t > 0$.

Following ([4], Section IV.C), a global effective envelope $\mathcal{H}_{\mathcal{C}}^{l, \varepsilon}$ for time intervals of length l can be constructed in the steps outlined below.

Step 1: Select a finite number of values for $\mathcal{H}_{\mathcal{C}}^{l, \varepsilon}$. We fix a set of values τ_1, \dots, τ_n so that $0 \leq \tau_1 < \dots < \tau_n = l$, and a set of integers k_1, \dots, k_n . We define

$$\mathcal{H}_{\mathcal{C}}^{l, \varepsilon}(\tau_i) = \mathcal{G}_{\mathcal{C}}^{\varepsilon}(\tau_i'), \quad (43)$$

where

$$\tau'_i = \frac{k_i + 1}{k_i} \tau_i \quad \text{and} \quad \varepsilon' = \varepsilon \left(\sum_{i=1}^n \frac{lk_i}{\tau_i} \right)^{-1}. \quad (44)$$

It was shown in [4] that for any interval I_l of length l ,

$$\Pr \left\{ A_{\mathcal{C}}(t, t + \tau_i) \leq \mathcal{H}_{\mathcal{C}}^{l, \varepsilon}(\tau_i), \forall t, \tau_i : [t, t + \tau_i] \subseteq I_l \right\} \geq 1 - \varepsilon. \quad (45)$$

This creates n points of a global effective envelope. In the next steps, we generate from these points a bound for all points.

Remark:

To obtain a ‘good’ global effective envelope, it is important to make good choices for n, τ_1, \dots, τ_n , and k_1, \dots, k_n . In our examples in Section 6, we select the points following a heuristic from ([4], Appendix II). The heuristic is motivated by appealing to the Central Limit Theorem and selects the parameters as follows.

If $\mathcal{G}_{\mathcal{C}}^{\varepsilon}$ is constructed for an aggregate of flows following Subsection 5.1, and if the individual flows are peak rate constrained with peak rates P_j and with average rate ρ_j for flow j , we set

$$\gamma = 1 + \frac{1}{k + 1}, \quad (46)$$

$$k = \left\lfloor z \left(z + \frac{\sum_{j \in \mathcal{C}} \rho_j}{\sqrt{\sum_{j \in \mathcal{C}} \rho_j (P_j - \rho_j)}} \right) \right\rfloor, \quad (47)$$

where $\lfloor x \rfloor$ is the largest integer smaller than x , and where z is chosen so that $1 - \Phi(z) = \varepsilon$.

With these values of k and γ , we set

$$k_i = k \quad \text{and} \quad \tau_i = \gamma^i \tau_o \quad i = 1, \dots, n, \quad (48)$$

where τ_o is a small number, and n is the smallest integer so that $\gamma^n \tau_o \geq l$.

Step 2: Interpolation. We use the n points generated in Step 1 to build a complete function. We use assumption (A1) and, if available, the deterministic bound $A_{\mathcal{C}}^*$ to interpolate between the values of $\mathcal{H}_{\mathcal{C}}^{l, \varepsilon}(\tau_i)$ as follows:

$$f(\tau) := \begin{cases} \min \left[A_{\mathcal{C}}^*(\tau), \mathcal{H}_{\mathcal{C}}^{l, \varepsilon}(\tau_1) \right], & \text{if } \tau \in [0, \tau_1], \\ \min \left[\mathcal{H}_{\mathcal{C}}^{l, \varepsilon}(\tau_{i-1}) + A_{\mathcal{C}}^*(\tau - \tau_{i-1}), \mathcal{H}_{\mathcal{C}}^{l, \varepsilon}(\tau_i) \right], & \text{if } \tau \in [\tau_{i-1}, \tau_i]. \end{cases} \quad (49)$$

Note that this definition gives a finite value for $f(\tau)$ even when $A_{\mathcal{C}}^*(t) = \infty$ for some $t > 0$. As shown in [4], f satisfies for any interval I_l of length l

$$\Pr \left\{ A_{\mathcal{C}}(t, t + \tau) \leq f(\tau), \forall t, \tau : [t, t + \tau] \subseteq I_l \right\} \geq 1 - \varepsilon. \quad (50)$$

Thus, f is with high probability a bound for all points. However, f is not a global effective envelope, since f may not be a subadditive function.

Step 3: Enforcing the Subadditivity of $\mathcal{H}_C^{l,\varepsilon}$. It was shown in [4] that the largest subadditive function below f , given by

$$\mathcal{H}_C^{l,\varepsilon}(\tau) = \inf_{\{\theta_j\} : \sum \theta_i = \tau} \sum_i f(\theta_i), \quad (51)$$

is a global effective envelope for the arrivals A_C . This completes the construction.

Remarks:

- Steps 1 and 2 of the construction of $\mathcal{H}_C^{l,\varepsilon}$ are computationally easy, being linear in the number n of values used in Step 1. Step 3, on the other hand, is at least quadratic in n . However, the function f constructed in Step 2 is a convenient upper bound for $\mathcal{H}_C^{l,\varepsilon}$, which is tight if the local effective envelope $\mathcal{G}_C^\varepsilon$ is close to being a subadditive function.
- In some scenarios, one may want to construct an effective envelope for an aggregate of flows from the effective envelopes of subsets of the aggregate. Suppose that $\mathcal{G}_{C_1}^{\varepsilon_1}$ and $\mathcal{G}_{C_2}^{\varepsilon_2}$ are local effective envelopes for two sets C_1 and C_2 , then $\mathcal{G}_{C_1 \cup C_2}^{\varepsilon_1 + \varepsilon_2} = \mathcal{G}_{C_1}^{\varepsilon_1} + \mathcal{G}_{C_2}^{\varepsilon_2}$ is a local effective envelope for $C_1 \cup C_2$. The same holds for global effective envelopes.

6 Evaluation

We now present numerical examples which demonstrate different applications of effective service curves, and evaluate the statistical multiplexing gain feasible with effective service curves from Section 4.

We assume that individual flows are regulated at the entrance to the network, using a peak rate limited leaky bucket with arrival envelope $A_j^*(\tau) = \min\{P_j\tau, \sigma_j + \rho_j\tau\}$ for flow j , where $P_j \geq \rho_j$ is the peak rate, ρ_j is the average rate, and σ_j is a burst size parameter. We consider two types of flows with parameters as given in the following table:

Type	Peak Rate P_j (Mbps)	Mean Rate ρ_j (Mbps)	Burst Size σ_j (bits)
Type 1	1.5	0.15	95400
Type 2	6.0	0.15	10345

The parameters are selected to be equal to those in [4, 14] and other studies.⁴

We assume that the arrivals satisfy assumptions (A1)–(A4), and we construct effective envelopes as shown in Section 5. We assume that capacities are allocated to aggregates of flows, in terms of deterministic service curves S_C and \overline{S}_C for a set C of flows. We will assume that service curves for the aggregate have a very simple constant-rate form, such as $S_C(t) = Nc t$ ($c > 0$), where c is referred to as ‘per-flow capacity’, and $N = |C|$ is the number of flows. For the construction of minimum effective service curves $\mathcal{S}_j^\varepsilon$, unless specifically stated otherwise, we use the most conservative bound from Corollary 1.4, i.e., $\mathcal{S}_j^\varepsilon = [S_C - \mathcal{H}_C^{l,\varepsilon}]_+$. This bound does not require a maximum service curve (as used in Theorem 3) and merely

⁴In this section, we use A_1^* and A_2^* to denote the arrival envelope of a Type-1 and a Type-2 flow, respectively.

requires us to calculate the global effective envelope $\mathcal{H}_C^{l,\varepsilon}$.⁵ For the length of the intervals, l , in the construction of $\mathcal{H}_C^{l,\varepsilon}$, we chose $l = 8$ seconds which is verified to be a deterministic bound on the busy period in the sense of Eqn. (35) for all examples. Note that l can be significantly reduced by using the probabilistic bounds on the busy period of Theorem 4.

Throughout this section, ‘effective service curve’ always refers to the minimum effective service curve. We compare the results obtained with effective service curves to the following non-statistical per-flow service provisioning schemes.

- A *peak rate* allocation, where each flow j has a service curve of $S_j(t) = P_j t$, provides an upper bound for the amount of resources reserved for a flow.
- An *average rate* allocation, where each flow j has a service curve of $S_j(t) = \rho_j t$, is a lower bound for the amount of resources reserved.
- A *deterministic* allocation delivers worst-case delay guarantees. The resources allocated to a flow are determined by the smallest (deterministic) constant-rate service curve $S_j(t) = \hat{c}_j t$ that satisfies the delay bound d , i.e., $\hat{c}_j = \inf \{c \geq 0 \mid \forall t \geq 0 : A^*(t - d) \leq c t\}$.

6.1 Example 1: Single Node

We investigate arrivals from a group \mathcal{C} of N flows with a delay bound of $d = 50$ ms at a single node. The flows are either all Type-1 or all Type-2 flows. We first compare the shape of effective service curves for different values of N and for $\varepsilon = 10^{-9}$, with the deterministic service curves.

We assume that the capacity allocated for the aggregate of N flows is $S_C(t) = N\hat{c}t$, where \hat{c} is the constant-rate service curve required by a flow to satisfy a delay bound of $d = 50$ ms according to the deterministic allocation given above. The required constant-rates are $\hat{c} \approx 0.8785$ for Type-1 flows, and $\hat{c} \approx 0.2050$ for Type-2 flows. Next, we calculate, via Theorem 3 and Corollary 1.4, effective service curves for any single flow from the set. That is, we calculate the effective curve as

1. $S_j^\varepsilon(t) = [N\hat{c}t - \mathcal{H}_{C-\{j\}}^{l,\varepsilon} * \overline{S}_{C-\{j\}}(t)]_+$ according to Theorem 3, or
2. $S_j^\varepsilon(t) = [N\hat{c}t - \mathcal{H}_C^{l,\varepsilon}(t)]_+$, according to Corollary 1.4.

Recall that Theorem 3 yields, among the effective service curves derived in Section 4, the least conservative and Corollary 1.4 the most conservative effective service curve.

In Figures 2(a) and 2(b) we plot effective service curves S_j^ε for different values of N , and compare them to the deterministic service curve $S_j(t) = \hat{c}t$. Effective service curves computed according to Theorem 3 are shown as dotted lines, effective service curves computed according to Corollary 1.4 are included as solid lines, and deterministic service curves $S_j(t) = \hat{c}t$ are included as dashed lines.

Figures 2(a) and 2(b) illustrate that, for large N , the effective service curves are significantly larger than service curves obtained from a deterministic allocation. The figures also show that the effective service curves generated from Corollary 1.4, at least for some values of N , are not monotonic as a function of τ . This can be explained by the fact that the global effective envelope is not concave, but only subadditive. Note that the difference between a subadditive function, i.e., $\mathcal{H}_C^{l,\varepsilon}$, and a constant rate function, i.e., $N\hat{c}t$, need

⁵In our examples, the numerical computations for effective envelopes are done in discrete intervals of length $\Delta = 1.3$ ms, and not in continuous time. This may introduce discretization errors.

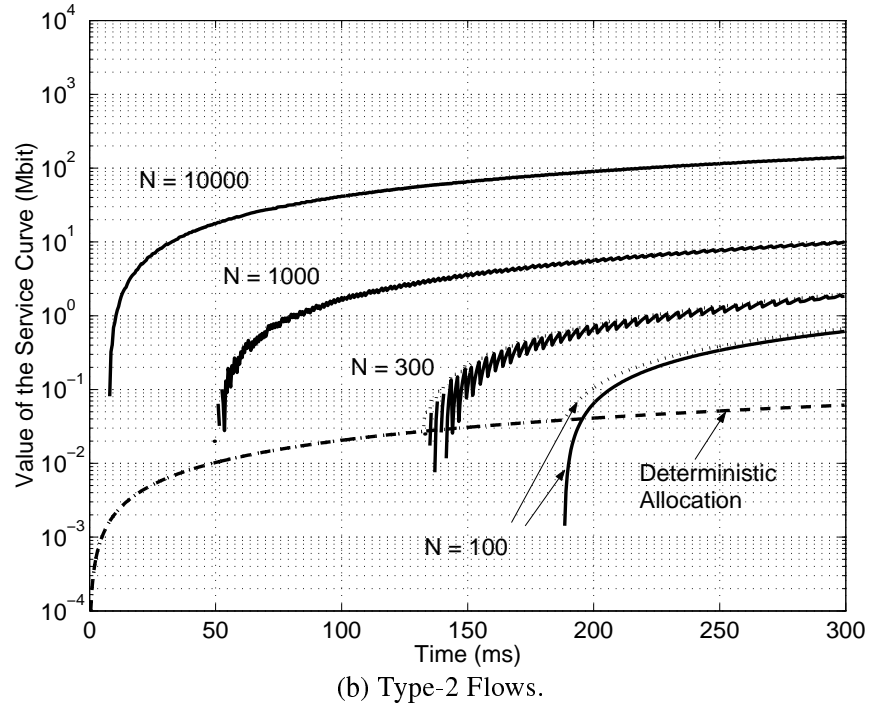
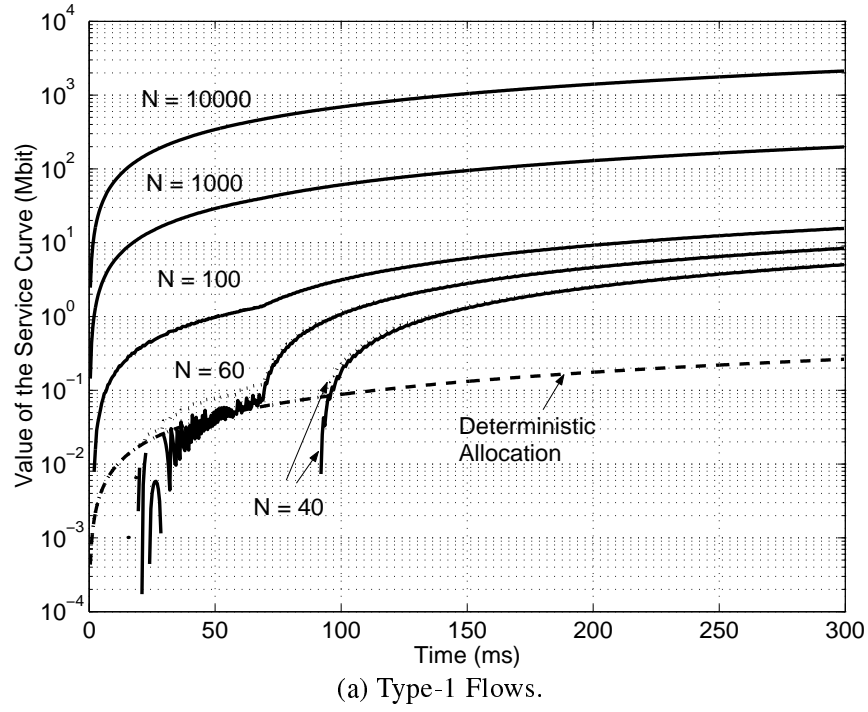


Figure 2: Example 1: Effective vs. deterministic service curve as a function of time, with capacity per flow computed for a deterministic delay bound of 50 ms. Effective service curves are shown for different values of N , and are calculated for $\varepsilon = 10^{-9}$.

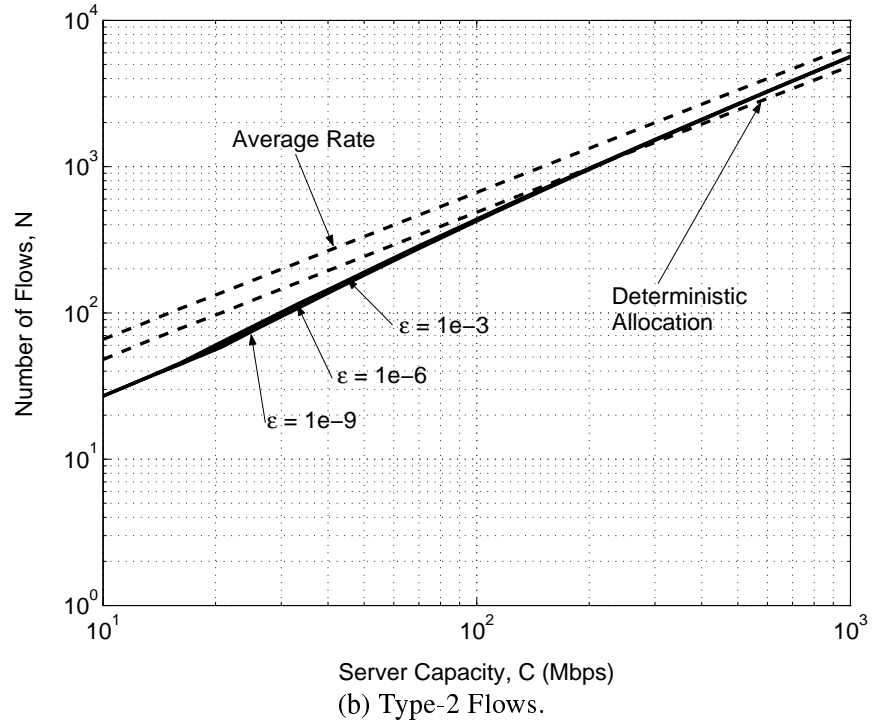
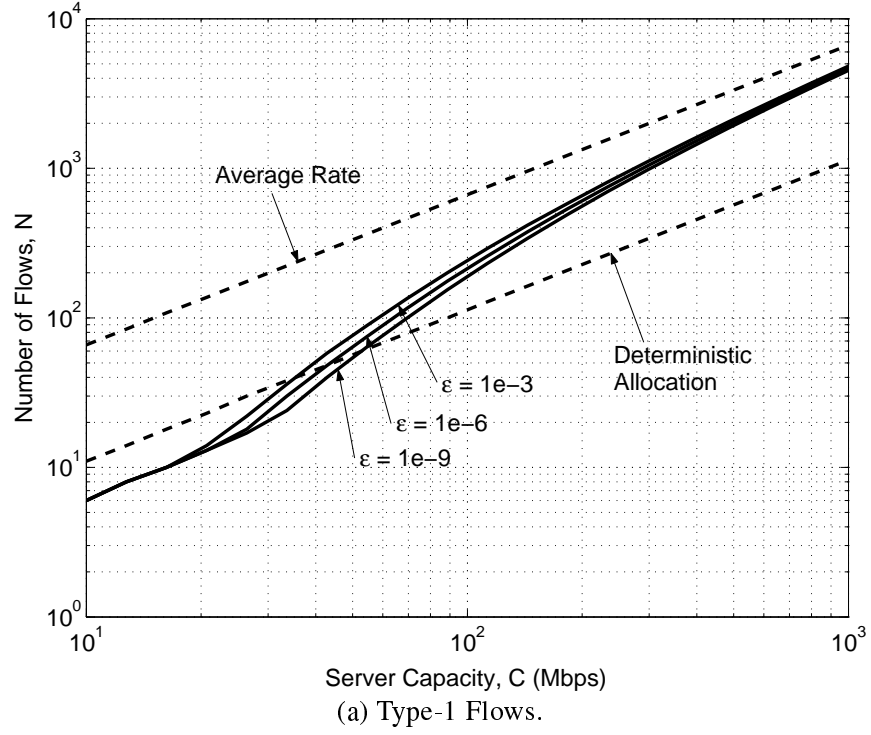


Figure 3: Example 1: Number of flows admitted on a link with capacity C to satisfy a delay bound of $d = 50$ ms with probability $1 - \epsilon$.

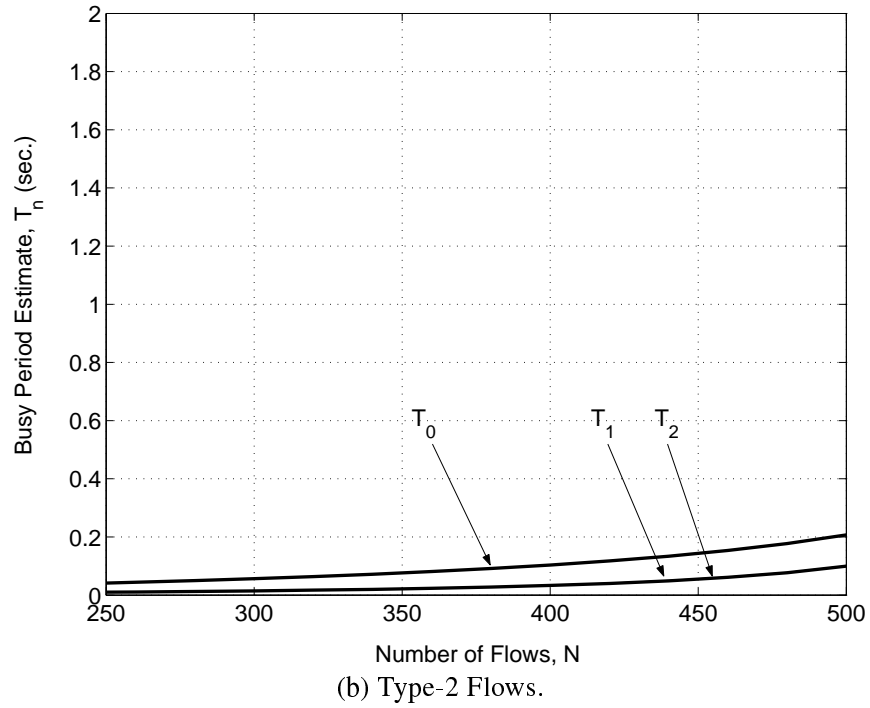
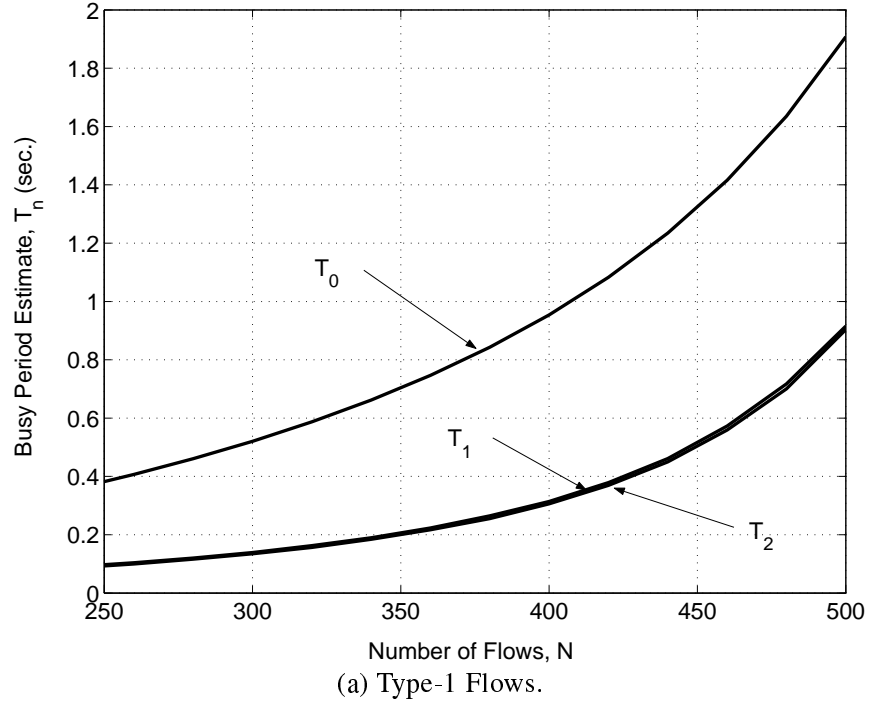


Figure 4: Example 1: Busy period estimates according to Theorem 4 on a link with a capacity $C = 100$ Mbps. The bound T_n is a bound for the busy period with probability $1 - n\varepsilon$.

not be monotonic. Here, it may be desirable to substitute the global effective envelopes in the construction of $\mathcal{S}_j^\varepsilon$ with concave upper bounds. To explain the visible change of slope of the curve for $N = 60$ and $N = 100$ at $t \approx 70$ ms in Figure 2(a), we note that at $\sigma/(P - \rho) \approx 70$ ms, the arrival envelope of Type-1 flows changes from Pt to $\sigma + \rho t$.

We note that the results from Theorem 3 and Corollary 1.4 can be distinguished only for small values of N . Since the effective service curves with Corollary 1.4 are easier to compute and do not require knowledge of a maximum service curve for the aggregate, from now on, we will consider only the most conservative effective service curves from Corollary 1.4.

In Figures 3(a) and 3(b), we evaluate the number of flows that can be provisioned on a link with capacity C , again using a delay bound of $d = 50$ ms. We use the same deterministic service curve as before, i.e., $S_j(t) \approx 0.8785t$ for Type-1 flows and $S_j(t) \approx 0.2050t$ for Type-2 flows. For the effective service curve we find the largest N such that $\mathcal{S}_j^\varepsilon(t) = [Ct - \mathcal{H}_C^{l,\varepsilon}(t)]_+$ assures via Theorem 2 the delay bound d with probability $1 - \varepsilon$ (recall that $N = |\mathcal{C}|$). The figures include plots of effective service curves with $\varepsilon = 10^{-3}, 10^{-6}, 10^{-9}$, along with the result for an average rate allocation (which does not satisfy the delay bound). The graphs show that, even for ε very small, the effective service curve shows significant statistical multiplexing gain, as C is increased. For large C , the plots for effective service curves and average rate allocation become close for both types of flows. For small C , on the other hand, the number of flows is too small to extract multiplexing gain, and, consequently, the effective service curve constructed from Corollary 1.4 can be inferior to a deterministic service curve.

Finally, we will illustrate the probabilistic bounds on the busy periods from Theorem 4. For a link with capacity $C = 100$ Mbps, we calculate the deterministic bound on the length of a busy period according to Eqn. (35), and compare this bound with the probabilistic bounds provided by Theorem 4. In Figure 4 we evaluate the bounds on busy periods for values of $N = 250 - 500$ Type-1 flows or Type-2 flows. Recall that T_n is a bound for the length of the busy period with probability $1 - n\varepsilon$. The plots indicate that T_1 significantly reduces the busy period bound, but that further iterations of Theorem 4 do not result in improvements.

6.2 Example 2: Multiple Nodes with Cross Traffic

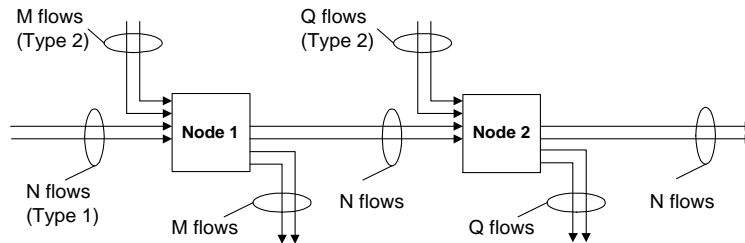
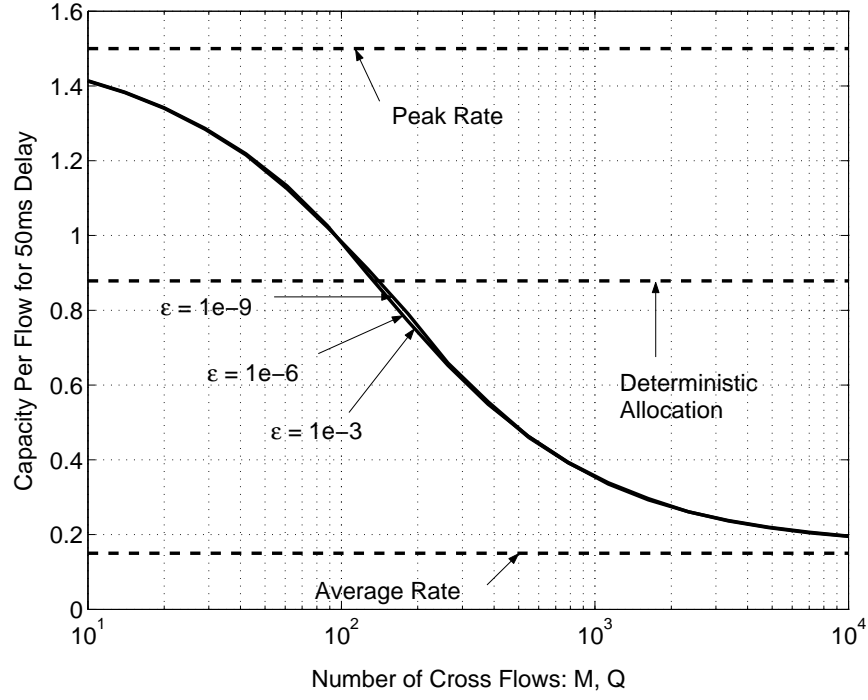
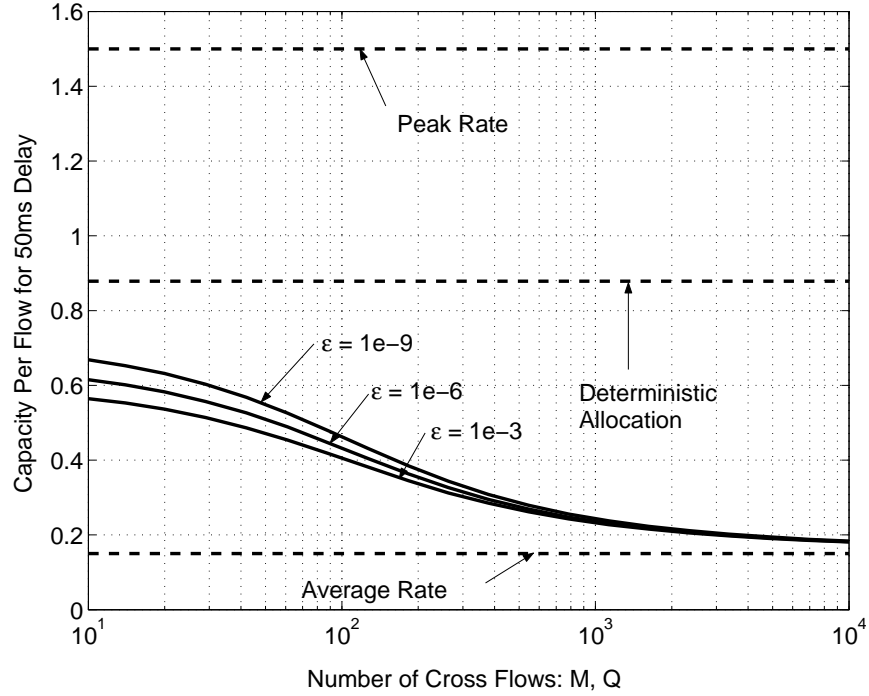


Figure 5: Example 2: A network with 2 nodes and with cross traffic.

We consider a network with two nodes, as shown in Figure 5, and determine the multiplexing gain attainable with effective service curves for a set of N flows through these two nodes, with an end-to-end delay bound of $d = 50$ ms. We assume there is cross traffic from M flows at the first node, and from Q flows at the second node. We assume that all N flows are Type-1 flows and that the cross traffic consists of Type-2 flows. We will denote flows that pass through both nodes as ‘through flows’ and cross traffic as ‘cross flows’.



(a) Approach 1: We use $\mathcal{H}_C^{2;l,\epsilon} = (N A_1^*) \oslash \mathcal{S}_C^{1,\epsilon}$.



(b) Approach 2: Through flows are rate-controlled at Node 2 with policing function $A_1^* \oslash \mathcal{S}_1^{1,\epsilon}$.

Figure 6: Example 2: Required capacity for 100 Type-1 flows to obtain an end-to-end delay bound of $d = 50$ ms with probability $1 - 2\epsilon$, as a function of the number of Type-2 cross flows. The plots show a comparison of two different approaches described in the text.

In this example, we make for convenience, a small change of notation. We let S_{N+M}^1 denote the deterministic service curve allocated to the aggregate of N Type-1 and M Type-2 flows at Node 1. Likewise, we let S_{N+Q}^2 denote the deterministic service curve allocated to the aggregate of N Type-1 and Q Type-2 flows at Node 2. $\mathcal{H}_{N+M}^{1;l,\varepsilon}$ will denote the global effective envelope of the $N + M$ flows at the first node, and $\mathcal{H}_{N+Q}^{2;l,\varepsilon}$ denotes the global effective envelope of the $N + Q$ flows at the second node.

Using Corollary 1.4, an effective service curve of a Type-1 flow j can be given by $S_j^{1,\varepsilon} = [S_{N+M}^1 - \mathcal{H}_{N+M}^{1;l,\varepsilon}]_+$ and $S_j^{2;l,\varepsilon} = [S_{N+Q}^2 - \mathcal{H}_{N+Q}^{2;l,\varepsilon}]_+$. We set $S_{N+Q}^1 = \tilde{c}(N + M)t$ and $S_{N+Q}^2 = \tilde{c}(N + Q)t$, where $\tilde{c} > 0$ is selected as the smallest value such that the effective network service curve of a Type-1 flow, $S_j^{net,2\varepsilon} = S_j^{1,\varepsilon} * S_j^{2;l,\varepsilon}$, satisfies a probabilistic end-to-end delay bound of $d = 50$ ms, according to Theorem 2. Recall, from Theorem 2 that the probability of an end-to-end delay bound violation is 2ε .

Note that there are several alternatives for calculating the global effective envelope $\mathcal{H}_{N+Q}^{2;l,\varepsilon}$ at the second node.

- **Approach 1:** We can determine $\mathcal{H}_{N+Q}^{2;l,\varepsilon}$ by adding the global effective envelopes of the N through flows and the Q cross flows. We have $\mathcal{H}_{N+Q}^{2;l,\varepsilon} = \mathcal{H}_N^{2;l,\varepsilon/2} + \mathcal{H}_Q^{2;l,\varepsilon/2}$, where $\mathcal{H}_N^{2;l,\varepsilon/2}$ and $\mathcal{H}_Q^{2;l,\varepsilon/2}$ are the global effective envelopes for the arrivals of the through flows and the cross flows, respectively, at the second node. Recall from the remark at the end of Section 5 that constructing an effective envelope from the addition of effective envelopes requires us to split the ε . Here, $\mathcal{H}_Q^{2;l,\varepsilon}$ is computed as described in Section 5. We can calculate an effective service curve for the aggregate of the N through flows at the first node as $S_N^{1,\varepsilon} = [S_{N+M}^1 - \mathcal{H}_M^{1;l,\varepsilon}]_+$, where $\mathcal{H}_M^{1;l,\varepsilon}$ is the global effective envelope of the M cross flows at the first node. Then we can obtain $\mathcal{H}_N^{2;l,\varepsilon}$ as follows. Per Theorem 2, the function $(N \cdot A_1^*) \odot S_N^{1,\varepsilon}$ yields a local effective envelope, which can be turned into a global effective envelope using Subsection 5.2. We choose the construction parameter k based on the peak and average rates of the arrivals to the first node.
- **Approach 2:** This approach requires a few additional assumptions. From Theorem 2, we have that $A_1^* \odot S_j^{1,\varepsilon}$ is a local effective envelope for the departures of a Type-1 flow from Node 1. Suppose we can police or shape each through flow before it enters Node 2, and discard or delay all traffic exceeding $A_1^* \odot S_j^{1,\varepsilon}$. If we further assume that traffic from through flows which arrives to Node 2 from their respective policer or shaper are independent in the sense of assumption (A4) from Section 5, then we can construct $\mathcal{H}_{N+Q}^{2;l,\varepsilon}$ for all $N + Q$ flows at the second node, using A_2^* as arrival envelope for each cross flow, and $A_1^* \odot S_j^{1,\varepsilon}$ as the arrival envelope for each through flow. Note that adding a policer introduces losses and adding a shaper introduces additional delays, which are not accounted for.

In Figure 6, we show the results for the two approaches for an example with $N = 100$ through flows. We plot the required per-flow capacity \tilde{c} as a function of M and Q for $\varepsilon = 10^{-3}, 10^{-6}, 10^{-9}$. The number of cross traffic flows is assumed to be identical at both nodes, that is, $M = Q$, and varied from 10 to 10,000 flows. We also consider results for a peak rate allocation with $P = 1.5$ Mbps, average rate allocation $\rho = 0.15$ Mbps, and a deterministic allocation $\hat{c} \approx 0.8785$ Mbps.

The figure illustrates that the required bandwidth to satisfy an end-to-end delay bound of $d = 50$ ms is close to an average rate allocation when the number of cross flows is large. In Figure 6(b), we see that Approach 2, which assumes shaping/policing of through flows at the second node yields a better multiplexing gain. The explanation for the difference is that calculating an effective envelope, here $\mathcal{H}_{N+Q}^{2;l,\varepsilon}$, by adding global effective envelopes of subsets of the aggregate, is less effective in exploiting multiplexing gain.

6.3 Example 3: Multiple Nodes With No Cross Traffic

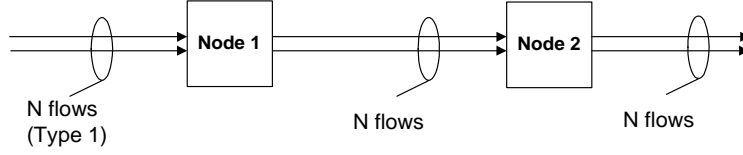


Figure 7: Example 3: A network with 2 nodes and no cross traffic.

We consider the two-node network shown in Figure 7 with no cross traffic, with N flows passing through both nodes. We will again evaluate the per-flow capacity needed at each node to satisfy a probabilistic or deterministic end-to-end delay bound of $d = 50$ ms. Similar to Example 2, we set $S_C^1(t) = S_C^2(t) = \tilde{c} N t$ to be the deterministic service curves allocated to the flows at the two nodes, where $\tilde{c} > 0$ is set to be the smallest rate such that the end-to-end delay bounds are satisfied.

Since there is no cross traffic, the analysis of this example can be much simplified by referring to the deterministic network calculus. For the aggregate of N flows, we obtain from Theorem 1 that $S_C^{net} = S_C^1 * S_C^2$ is a deterministic effective service curve. Since $S_C^1 = S_C^1 * S_C^2$, we obtain $S_C^{net} = \tilde{c} N t$ as deterministic network service curve for all flows. With this, we can now proceed as in Example 1, that is, construct an effective service curve with $S_j^{net, \varepsilon}(t) = [\tilde{c} N t - \mathcal{H}_C^{1:l, \varepsilon}(t)]_+$.

We assume that the N flows are either all Type-1 flows or all Type-2 flows. We graph the per-flow capacity \tilde{c} required to satisfy a probabilistic delay bound, $d = 50$ ms, as a function of N . Note that, with the above argument, the probability for meeting the end-to-end delay bound is $1 - \varepsilon$. The results are shown in Figures 8(a) and (b) for Type-1 and Type-2 flows, respectively. The plots illustrate that the bandwidth requirements of a flow approach the average rate, as the number of flows is increased. Interestingly, Figure 8(a) shows that the capacity requirement using effective service curves can exceed the peak rate for very small values of N .

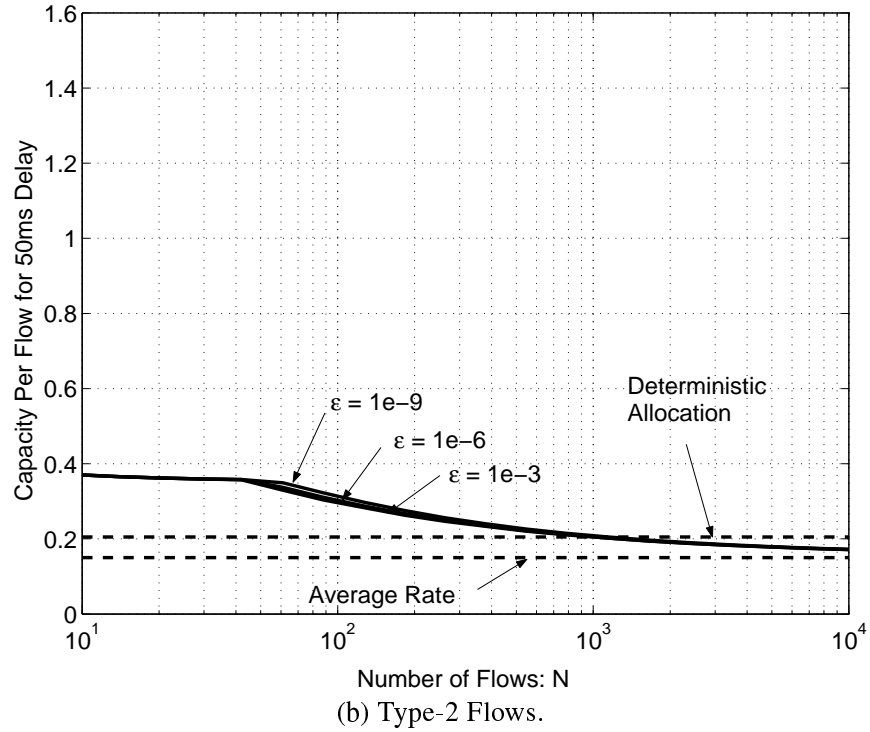
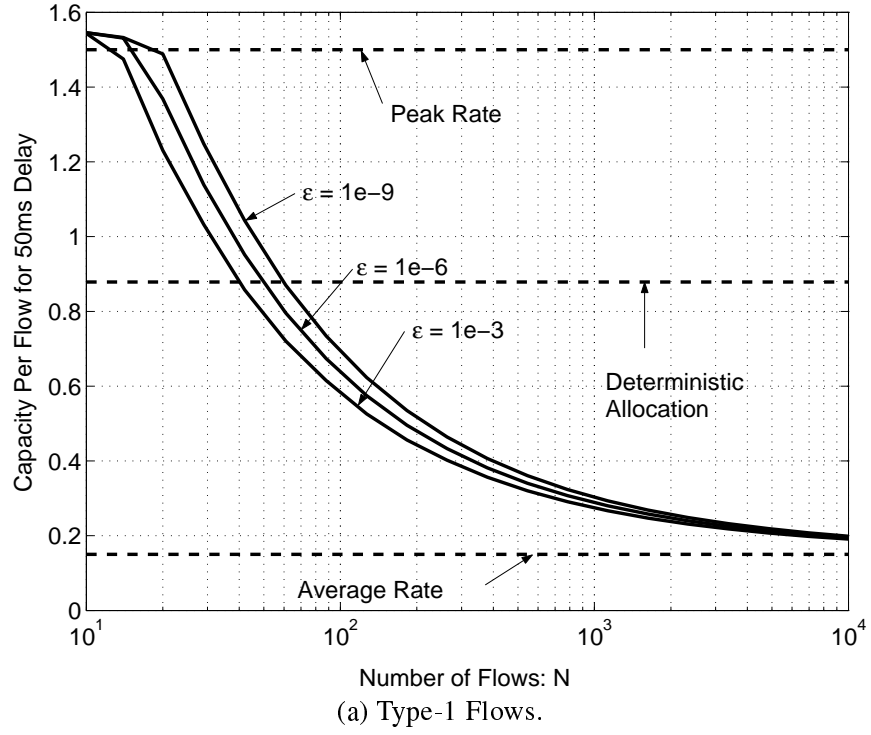


Figure 8: Example 3: Capacity per flow needed to support a delay bound of $d = 50$ ms with probability $1 - \varepsilon$ as a function of the total number of Type-1 flows. The effective service curve is calculated as $S_j^{net, \varepsilon}(t) = [\tilde{c} N t - \mathcal{H}_C^{1;l, \varepsilon}(t)]_+$.

7 Conclusions

We have presented a network calculus for statistically multiplexed traffic, which introduces the notion of *effective service curves* as a probabilistic bound on the service received by individual flows in a network. We have shown that some key results from the deterministic network calculus can be carried over to the statistical framework by inserting appropriate probabilistic arguments. Through the use of effective service curves, we are able to describe the service delivered to individual flows when capacities are allocated to aggregates of flows.

This paper raises a number of directions for future research:

- In addition to the statistical network calculus presented here, one can derive probabilistic bounds on properties of aggregates of flows, such as bounds on the output from a node, the backlog, and the delay at a node. In this calculus, given a set \mathcal{C} of flows, and a global envelope $\mathcal{H}_{\mathcal{C}}^{l,\varepsilon}$ for the arrivals from \mathcal{C} , we assume deterministic minimum (maximum) service curves $S_{\mathcal{C}}$ ($\overline{S}_{\mathcal{C}}$) which gives deterministic bounds on the service allocated to the aggregate of the flows in \mathcal{C} . It can be shown that (1) $\mathcal{H}_{\mathcal{C}}^{l,\varepsilon} \otimes S_{\mathcal{C}}$ is a global effective envelope for the output, (2) $b_{max} = \mathcal{H}_{\mathcal{C}}^{l,\varepsilon} \otimes S_{\mathcal{C}}(0)$ is a probabilistic backlog bound in $[0, l]$, and (3) $d_{max} = \inf \left\{ d \geq 0 \mid \forall 0 \leq t \leq l : \mathcal{H}_{\mathcal{C}}^{l,\varepsilon}(t - d) \leq S_{\mathcal{C}}(t) \right\}$ is a probabilistic delay bound for the interval $[0, l]$.
- Similar to global effective envelopes which are arrival bounds over all subintervals of an interval of length l , one can define global effective service curves, which give lower or upper bounds on the service received in all subintervals of a larger interval. In fact, the discussion at the beginning of Subsection 4.3 and Eqn. (32), indicate that a modified effective service curve similar to that in Theorem 3 may satisfy ‘global’ properties.
- The calculus in Section 3 works with a deterministic arrival bound and a probabilistic service curve. It may be of interest to study a statistical network calculus where both arrivals and service curves are viewed in terms of stochastic processes.
- The calculus in this paper may be sufficient to provision end-to-end delays in feedforward networks, however, it is not directly applicable to networks with arbitrary topology and arbitrary routes. The problems encountered should be similar to those in the deterministic network calculus [10, 18].

References

- [1] R. Agrawal, R. L. Cruz, C. Okino, and R. Rajan. Performance bounds for flow control protocols. *IEEE/ACM Transactions on Networking*, 7(3):310–323, June 1999.
- [2] M. Andrews. Probabilistic end-to-end delay bounds for earliest deadline first scheduling. In *Proceedings of IEEE Infocom 2000*, pages 603–612, Tel Aviv, March 2000.
- [3] F. L. Baccelli, G. Cohen, G. J. Olsder, and J.-P. Quadrat. *Synchronization and Linearity: An Algebra for Discrete Event Systems*. John Wiley and Sons, 1992.
- [4] R. Boorstyn, A. Burchard, J. Liebeherr, and C. Oottamakorn. Statistical service assurances for traffic scheduling algorithms. *IEEE Journal on Selected Areas in Communications. Special Issue on Internet QoS*, 18(12):2651–2664, December 2000.

- [5] J. Y. Le Boudec. Application of network calculus to guaranteed service networks. *IEEE/ACM Transactions on Information Theory*, 44(3):1087–1097, May 1998.
- [6] J. Y. Le Boudec and P. Thiran. *Network Calculus*. Springer Verlag, Lecture Notes in Computer Science, LNCS 2050, 2001.
- [7] C. S. Chang. Stability, queue length, and delay of deterministic and stochastic queueing networks. *IEEE Transactions on Automatic Control*, 39(5):913–931, May 1994.
- [8] C. S. Chang. On deterministic traffic regulation and service guarantees: a systematic approach by filtering. *IEEE/ACM Transactions on Information Theory*, 44(3):1097–1110, May 1998.
- [9] C. S. Chang. *Performance Guarantees in Communication Networks*. Springer Verlag, 2000.
- [10] R. L. Cruz. A Calculus for Network Delay, Part II: Network Analysis. *IEEE Transactions on Information Theory*, 37(1):132–141, January 1991.
- [11] R. L. Cruz. Quality of service guarantees in virtual circuit switched networks. *IEEE Journal on Selected Areas in Communications*, 13(6):1048–1056, August 1995.
- [12] R. L. Cruz. Quality of service management in integrated services networks. In *Proceedings of the 1st Semi-Annual Research Review, CWC, UCSD*, June 1996.
- [13] A. Elwalid and D. Mitra. Design of generalized processor sharing schedulers which statistically multiplex heterogeneous QoS classes. In *Proceedings of IEEE INFOCOM'99*, pages 1220–1230, New York, March 1999.
- [14] A. Elwalid, D. Mitra, and R. Wentworth. A new approach for allocating buffers and bandwidth to heterogeneous, regulated traffic in an ATM node. *IEEE Journal on Selected Areas in Communications*, 13(6):1115–1127, August 1995.
- [15] G. Kesidis and T. Konstantopoulos. Extremal traffic and worst-case performance for queues with shaped arrivals. In *Proceedings of Workshop on Analysis and Simulation of Communication Networks*, Toronto, November 1998.
- [16] J. Kurose. On computing per-session performance bounds in high-speed multi-hop computer networks. In *ACM Sigmetrics'92*, pages 128–139, 1992.
- [17] C. Li and E. Knightly. Coordinated network scheduling: A framework for end-to-end services. In *Proceedings of IEEE ICNP 2000*, Osaka, November 2000.
- [18] A. K. Parekh and R. G. Gallager. A generalized processor sharing approach to flow control in integrated services networks: The multiple node case. *IEEE/ACM Transactions on Networking*, 2(2):137–150, April 1994.
- [19] J. Qiu and E. Knightly. Inter-class resource sharing using statistical service envelopes. In *Proceedings of IEEE Infocom '99*, pages 36–42, March 1999.
- [20] V. Sivaraman and F. M. Chiussi. Statistical analysis of delay bound violations at an earliest deadline first scheduler. *Performance Evaluation*, 36(1):457–470, 1999.

- [21] V. Sivaraman and F. M. Chiussi. Providing end-to-end statistical delay guarantees with earliest deadline first scheduling and per-hop traffic shaping. In *Proceedings of IEEE Infocom 2000*, pages 603–612, Tel Aviv, March 2000.
- [22] D. Starobinski and M. Sidi. Stochastically bounded burstiness for communication networks. In *Proceedings of IEEE Infocom '99*, pages 36–42, March 1999.
- [23] O. Yaron and M. Sidi. Performance and stability of communication networks via robust exponential bounds. *IEEE/ACM Transactions on Networking*, 1(3):372–385, June 1993.

APPENDIX

The following sections include the proofs of Theorems 1, 2, and 3. Note that Theorem 1 states well-known results from the deterministic network calculus [1, 5, 8]. We follow the derivations in [1].

As in [1], we will use in the proofs of the delay bound the impulse function δ_τ , defined as

$$\delta_\tau(t) = \begin{cases} \infty, & \text{if } t > \tau, \\ 0, & \text{if } t \leq \tau. \end{cases} \quad (52)$$

Note that $f(t - \tau) = f * \delta_\tau(t)$.

A Proof of Theorem 1

In this section we present the proof of Theorem 1, following [1]. We emphasize that the proof of Theorem 2 uses the same arguments as the proofs of the deterministic network calculus.

A.1 Theorem 1: Proof of Output Bound

The departures $D(t, t + \tau)$ in the interval $[t, t + \tau)$ can be bounded, for all $t \geq 0$ and $\tau > 0$ as follows.

$$D(t, t + \tau) = D(t + \tau) - D(t) \quad (53)$$

$$\leq D(t + \tau) - A * S(t) \quad (54)$$

$$\leq A(t + \tau) - A * S(t) \quad (55)$$

$$= A(t + \tau) - \inf_{x \in [0, t]} [A(t - x) + S(x)] \quad (56)$$

$$= \sup_{x \in [0, t]} [A(t + \tau) - A(t - x) - S(x)] \quad (57)$$

$$= \sup_{x \in [0, t]} [A(t - x, t + \tau) - S(x)] \quad (58)$$

$$\leq \sup_{x \geq 0} [A(t - x, t + \tau) - S(x)] \quad (59)$$

$$\leq \sup_{x \geq 0} [A^*(\tau + x) - S(x)] \quad (60)$$

$$= A^* \oslash S(\tau). \quad (61)$$

Eqn. (54) follows from the definition of the minimum service curve S . Eqn. (55) uses that departures in $[0, t)$ cannot exceed arrivals, that is, $D(t) \leq A(t)$ for all $t \geq 0$. Eqn. (56) expands the convolution operator. Eqn. (57) takes $A(t + \tau)$ inside the infimum. Eqn. (58) uses that $A(t - x, t + \tau) = A(t + \tau) - A(t - x)$. Eqn. (59) extends the range of the supremum. Eqn. (60) uses the definition of an arrival envelope. Finally, Eqn. (61) uses the definition of the deconvolution operator.

A.2 Theorem 1: Proof of Backlog Bound

Recall that $B(t)$ denotes the backlog at time $t \geq 0$. The following establishes a bound for $B(t)$, for all $t \geq 0$ and $\tau > 0$.

$$B(t) = A(t) - D(t) \quad (62)$$

$$\leq A(t) - A * S(t) \quad (63)$$

$$= A(t) - \inf_{x \in [0, t)} [A(t - x) + S(x)] \quad (64)$$

$$= \sup_{x \in [0, t)} [A(t) - A(t - x) - S(x)] \quad (65)$$

$$= \sup_{x \in [0, t)} [A(t - x, t) - S(x)] \quad (66)$$

$$\leq \sup_{x \geq 0} [A(t - x, t) - S(x)] \quad (67)$$

$$\leq \sup_{x \geq 0} [A^*(x) - S(x)] \quad (68)$$

$$= A^* \oslash S(0) . \quad (69)$$

Eqn. (63) uses the definition of the service curve S . Eqn. (64) expands the convolution operator. Eqn. (65) takes $A(t)$ inside the infimum. Eqn. (66) applies the notation $A(t - x, t) = A(t) - A(t - x)$. Eqn. (67) extends the range of the supremum. Eqn. (68) uses the definition of an arrival envelope. Finally, Eqn. (69) uses the definition of the deconvolution operator.

A.3 Theorem 1: Proof of Delay Bound

The delay bound is proven by showing that for each $t \geq 0$, the output at time t , $D(t)$, is larger than the arrivals until time $t - d_{max}$, $A(t - d_{max})$. Thus, no backlogged traffic has violated the delay bound at time t .

$$D(t) \geq A * S(t) \quad (70)$$

$$\geq A * (A^* * \delta_{d_{max}})(t) \quad (71)$$

$$= (A * A^*) * \delta_{d_{max}}(t) \quad (72)$$

$$\geq A * \delta_{d_{max}}(t) \quad (73)$$

$$= A(t - d_{max}) . \quad (74)$$

Eqn. (70) uses the definition of the service curve S . Eqn. (71) follows from the given definition of d_{max} in Theorem 1, using the identity $f(t - \tau) = f * \delta_\tau(t)$. Eqn. (72) holds due to the associativity property of the convolution operator [3]. Eqn. (73) follows directly from the definition of the arrival envelope, $A(t) \leq A * A^*(t)$. Eqn. (74) uses the identity $f(t - \tau) = f * \delta_\tau(t)$.

A.4 Theorem 1: Proof of Network Service Curve

We first consider the case where the path through the network consists of just two nodes, $H = 2$. We will show that the network as a whole delivers the minimum service curve $S^{net} = S^1 * S^2$, such that for all

$t \geq 0$,

$$D^{net}(t) \geq A^{net} * (S^1 * S^2)(t) . \quad (75)$$

The proof for maximum service curves is analogous, and will be omitted.

Since the second node delivers the minimum service S^2 , we have for all $t > 0$.

$$D^{net}(t) \geq \inf_{x \in [0, t)} [A^2(x) + S^2(t - x)] \quad (76)$$

Also, since the first node delivers the minimum service S^1 , we have for all $x > 0$,

$$D^1(x) \geq \inf_{y \in [0, x)} [A^{net}(y) + S^1(x - y)] . \quad (77)$$

Since $A^2 = D^1$, we obtain from inserting Eqn. (77) into Eqn. (76) that for all $t \geq 0$,

$$D^{net}(t) \geq \inf_{x \in [0, t)} \left[\inf_{y \in [0, x)} [A^{net}(y) + S^1(x - y)] + S^2(t - x) \right] \quad (78)$$

In other words, for all $t \geq 0$,

$$D^{net}(t) \geq (A^{net} * S^1) * S^2(t) . \quad (79)$$

Finally, the associativity property of convolution [3] implies that for all $t \geq 0$,

$$D^{net}(t) \geq A^{net} * (S^1 * S^2)(t) . \quad (80)$$

Then, the claim that

$$D^{net}(t) \geq A^{net} * (S^1 * \dots * S^H)(t) . \quad (81)$$

for all $t \geq 0$ follows by induction.

B Proof of Theorem 2

B.1 Theorem 2: Proof of Output Bound

The derivations for the output bound start with $D(t, t + \tau) = D(t + \tau) - D(t)$ and the definition of the minimum effective service curve. For any fixed t and τ , we have

$$1 - \varepsilon \leq \Pr \{D(t, t + \tau) \leq D(t + \tau) - A * \mathcal{S}^\varepsilon(t)\} \quad (82)$$

$$\leq \Pr \{D(t, t + \tau) \leq A(t + \tau) - A * \mathcal{S}^\varepsilon(t)\} \quad (83)$$

$$= \Pr \left\{ D(t, t + \tau) \leq A(t + \tau) - \inf_{x \in [0, t]} [A(t - x) + \mathcal{S}^\varepsilon(x)] \right\} \quad (84)$$

$$= \Pr \left\{ D(t, t + \tau) \leq \sup_{x \in [0, t]} [A(t + \tau) - A(t - x) - \mathcal{S}^\varepsilon(x)] \right\} \quad (85)$$

$$= \Pr \left\{ D(t, t + \tau) \leq \sup_{x \in [0, t]} [A(t - x, t + \tau) - \mathcal{S}^\varepsilon(x)] \right\} \quad (86)$$

$$\leq \Pr \left\{ D(t, t + \tau) \leq \sup_{x \geq 0} [A(t - x, t + \tau) - \mathcal{S}^\varepsilon(x)] \right\} \quad (87)$$

$$\leq \Pr \left\{ D(t, t + \tau) \leq \sup_{x \geq 0} [A^*(\tau + x) - \mathcal{S}^\varepsilon(x)] \right\} \quad (88)$$

$$= \Pr \{D(t, t + \tau) \leq A^* \oslash \mathcal{S}^\varepsilon(\tau)\} . \quad (89)$$

Eqn. (82) follows from the definition of the minimum service curve \mathcal{S}^ε . Eqn. (83) uses that departures in $[0, t)$ cannot exceed arrivals, that is, $D(t) \leq A(t)$ for all $t \geq 0$. Eqn. (84) expands the convolution operator. Eqn. (85) merely takes $A(t + \tau)$ inside the infimum. Eqn. (86) uses the notation $A(t - x, t + \tau) = A(t + \tau) - A(t - x)$. Eqn. (87) extends the range of the supremum. Eqn. (88) uses the definition of an arrival envelope. Finally, Eqn. (89) uses the definition of the deconvolution operator.

B.2 Theorem 2: Proof of Backlog Bound

Since $B(t) = A(t) - D(t)$ and with the definition of the minimum effective service curve, we can write

$$1 - \varepsilon \leq \Pr \{B(t) \leq A(t) - A * \mathcal{S}^\varepsilon(t)\} \quad (90)$$

$$= \Pr \left\{ B(t) \leq A(t) - \inf_{x \in [0, t]} [A(t - x) + \mathcal{S}^\varepsilon(x)] \right\} \quad (91)$$

$$= \Pr \left\{ B(t) \leq \sup_{x \in [0, t]} [A(t) - A(t - x) - \mathcal{S}^\varepsilon(x)] \right\} \quad (92)$$

$$= \Pr \left\{ B(t) \leq \sup_{x \in [0, t]} [A(t - x, t) - \mathcal{S}^\varepsilon(x)] \right\} \quad (93)$$

$$\geq \Pr \left\{ B(t) \leq \sup_{x \geq 0} [A(t - x, t) - \mathcal{S}^\varepsilon(x)] \right\} \quad (94)$$

$$\leq \Pr \left\{ B(t) \leq \sup_{x \geq 0} [A^*(x) - \mathcal{S}^\varepsilon(x)] \right\} \quad (95)$$

$$= \Pr \{B(t) \leq A^* \oslash \mathcal{S}^\varepsilon(0)\} . \quad (96)$$

Eqn. (90) uses the definition of the minimum effective service curve \mathcal{S}^ε . Eqn. (91) expands the convolution operator. Eqn. (92) takes $A(t)$ inside the infimum. Eqn. (93) applies the notation $A(t - x, t) = A(t) - A(t - x)$. Eqn. (94) extends the range of the supremum. Eqn. (95) uses the definition of an arrival envelope. Finally, Eqn. (96) uses the definition of the deconvolution operator.

B.3 Theorem 2: Proof of Delay Bound

The delay bound is proven by showing that for any $t \geq 0$, the output at time t , $D(t)$, is larger than the arrivals until time $t - d_{max}$, $A(t - d_{max})$. Hence, no backlogged traffic at time t has violated the delay bound.

$$1 - \varepsilon \leq Pr \{D(t) \geq A * \mathcal{S}^\varepsilon(t)\} \quad (97)$$

$$\leq Pr \{D(t) \geq A * (A^* * \delta_{d_{max}})(t)\} \quad (98)$$

$$= Pr \{D(t) \geq (A * A^*) * \delta_{d_{max}}(t)\} \quad (99)$$

$$\leq Pr \{D(t) \geq A * \delta_{d_{max}}(t)\} \quad (100)$$

$$= Pr \{D(t) \geq A(t - d_{max})\} . \quad (101)$$

Eqn. (97) uses the definition of the minimum effective service curve \mathcal{S}^ε . Eqn. (98) follows from the definition of d_{max} from Theorem 2, and uses the identity $f(t - \tau) = f * \delta_\tau(t)$. Eqn. (99) holds due to the associativity property of the convolution operator [3]. Eqn. (100) follows from $A(t) \leq A * A^*(t)$ for all $t \geq 0$. Eqn. (101) again uses the identity $f(t - \tau) = f * \delta_\tau(t)$.

B.4 Theorem 2: Proof of Network Effective Service Curve

We first consider the case where the path through the network consists of just two nodes, $H = 2$. We show that the network as a whole delivers the minimum effective service curve $\mathcal{S}^{net, \varepsilon_1 + \varepsilon_2} = \mathcal{S}^{1, \varepsilon_1} * \mathcal{S}^{2, \varepsilon_2}$. Note that the violation probability for the network effective service curve will be $\varepsilon_1 + \varepsilon_2$, i.e.,

$$Pr \{D^{net}(t) \geq A^{net} * (\mathcal{S}^{1, \varepsilon_1} * \mathcal{S}^{2, \varepsilon_2})(t)\} \geq 1 - (\varepsilon_1 + \varepsilon_2). \quad (102)$$

The proof for maximum effective service curves is analogous and omitted.

Fix a value of $t > 0$. Since the second node delivers the minimum effective service $\mathcal{S}_2^{\varepsilon_2}$, we have

$$Pr \left\{ D^{net}(t) \geq \inf_{x \in [0, t]} \left[A^2(x) + \mathcal{S}^{2, \varepsilon_2}(t - x) \right] \right\} \geq 1 - \varepsilon_2. \quad (103)$$

Also, since the first node delivers the minimum effective service $\mathcal{S}^{1, \varepsilon_1}$, we have

$$Pr \left\{ D^1(x) \geq \inf_{y \in [0, x]} \left[A^{net}(y) + \mathcal{S}^{1, \varepsilon_1}(x - y) \right] \right\} \geq 1 - \varepsilon_1. \quad (104)$$

Therefore, the probability of a violation at any of the two nodes is bounded as follows.

$$Pr \left\{ \begin{array}{c} D^{net}(t) < \inf_{x \in [0, t]} \left[A^2(x) + \mathcal{S}^{2, \varepsilon_2}(t - x) \right] \\ \text{or} \\ D^1(x) < \inf_{y \in [0, x]} \left[A^{net}(y) + \mathcal{S}^{1, \varepsilon_1}(x - y) \right] \end{array} \right\} < \varepsilon_1 + \varepsilon_2, \quad (105)$$

Then, using the fact that $A^2 = D^1$, we obtain

$$Pr \left\{ \begin{array}{c} D^{net}(t) \geq \inf_{x \in [0,t]} [D^1(x) + \mathcal{S}^{2,\varepsilon_2}(t-x)] \\ \text{and} \\ D^1(x) \geq \inf_{y \in [0,x]} [A^{net}(y) + \mathcal{S}^{1,\varepsilon_1}(x-y)] \end{array} \right\} \geq 1 - (\varepsilon_1 + \varepsilon_2). \quad (106)$$

By inserting D^1 into the bound for D^{net} , we obtain

$$Pr \left\{ D^{net}(t) \geq \inf_{x \in [0,t]} \left[\inf_{y \in [0,x]} [A^{net}(y) + \mathcal{S}^{1,\varepsilon_1}(x-y)] + \mathcal{S}^{2,\varepsilon_2}(t-x) \right] \right\} \geq 1 - (\varepsilon_1 + \varepsilon_2). \quad (107)$$

We can rewrite the last equation as

$$Pr \{ D^{net}(t) \geq (A^{net} * \mathcal{S}^{1,\varepsilon_1}) * \mathcal{S}^{2,\varepsilon_2}(t) \} \geq 1 - (\varepsilon_1 + \varepsilon_2). \quad (108)$$

Finally, the associativity property of convolution [3] implies

$$Pr \{ D^{net}(t) \geq A^{net} * (\mathcal{S}^{1,\varepsilon_1} * \mathcal{S}^{2,\varepsilon_2})(t) \} \geq 1 - (\varepsilon_1 + \varepsilon_2). \quad (109)$$

Then, the claim

$$Pr \{ D^{net}(t) \geq A^{net} * (\mathcal{S}^{1,\varepsilon_1} * \dots * \mathcal{S}^{H,\varepsilon_H})(t) \} \geq 1 - (\varepsilon_1 + \dots \varepsilon_H). \quad (110)$$

follows by induction.

C Proof of Theorem 3

We will show that $\mathcal{S}_j^\varepsilon$ in Theorem 3 satisfies Definition 1.

$$D_j(t) = D_C - D_{C-\{j\}}(t) \quad (111)$$

$$\geq A_C * S_C(t) - A_{C-\{j\}} * \overline{S}_{C-\{j\}}(t) \quad (112)$$

$$= \inf_{x \in [0,t]} [A_C(t-x) + S_C(x)] - \inf_{y \in [0,t]} [A_{C-\{j\}}(t-y) + \overline{S}_{C-\{j\}}(y)] \quad (113)$$

$$= \inf_{x \in [0,t]} \left[A_C(t-x) + S_C(x) - \inf_{y \in [0,t]} [A_{C-\{j\}}(t-y) + \overline{S}_{C-\{j\}}(y)] \right] \quad (114)$$

$$\geq \inf_{x \in [0,t]} \left[A_C(t-x) + S_C(x) - \inf_{y \in [0,x]} [A_{C-\{j\}}(t-y) + \overline{S}_{C-\{j\}}(y)] \right] \quad (115)$$

$$= \inf_{x \in [0,t]} \left[(A_j(t-x) + A_{C-\{j\}}(t-x)) + S_C(x) - \inf_{y \in [0,x]} [(A_{C-\{j\}}(t-x) + A_{C-\{j\}}(t-x, t-y)) + \overline{S}_{C-\{j\}}(y)] \right] \quad (116)$$

$$= \inf_{x \in [0,t]} \left[A_j(t-x) + S_C(x) - \inf_{y \in [0,x]} [A_{C-\{j\}}(t-x, t-y) + \overline{S}_{C-\{j\}}(y)] \right], \quad (117)$$

where Eqn. (112) follows from the definition of minimum and maximum service curves, and Eqn. (113) merely expands the operators. Eqn. (114) rearranges the two infima without changing the range. Eqn. (115)

holds since the range of the second infimum is reduced, which may increase the infimum. Eqn. (116) uses the identities

$$\begin{aligned} A_{\mathcal{C}}(t-x) &= A_j(t-x) + A_{\mathcal{C}-\{j\}}(t-x) \\ A_{\mathcal{C}-\{j\}}(t-y) &= A_{\mathcal{C}-\{j\}}(t-x) + A_{\mathcal{C}-\{j\}}(t-x, t-y). \end{aligned}$$

Now we use that $\mathcal{H}_{\mathcal{C}-\{j\}}^{l,\varepsilon}$ is a global effective envelope for the flows in $\mathcal{C} - \{j\}$ for the interval $[0, l]$. Then, we have from Definition 2 that for

$$Pr \left\{ A_{\mathcal{C}-\{j\}}(\tau_1, \tau_2) \leq \mathcal{H}_{\mathcal{C}-\{j\}}^{l,\varepsilon}(\tau_2 - \tau_1), \forall \tau_1, \tau_2 : [\tau_1, \tau_2] \in [0, l] \right\} \geq 1 - \varepsilon. \quad (118)$$

Thus, we obtain from Eqn. (117) that

$$\begin{aligned} 1 - \varepsilon &\leq Pr \left\{ D_j(t) \geq \inf_{x \in [0, t)} \left[A_j(t-x) + S_{\mathcal{C}}(x) \right. \right. \\ &\quad \left. \left. - \inf_{y \in [0, x)} \left[\mathcal{H}_{\mathcal{C}-\{j\}}^{l,\varepsilon}(y-x) + \overline{S}_{\mathcal{C}-\{j\}}(y) \right] \right], \forall t : t \in [0, l] \right\} \end{aligned} \quad (119)$$

$$= Pr \left\{ D_j(t) \geq \inf_{x \in [0, t)} \left[A_j(t-x) + \left(S_{\mathcal{C}} - \mathcal{H}_{\mathcal{C}-\{j\}}^{l,\varepsilon} * \overline{S}_{\mathcal{C}-\{j\}} \right)(x) \right], \forall t : t \in [0, l] \right\} \quad (120)$$

$$= Pr \left\{ D_j(t) \geq A_j * \left(S_{\mathcal{C}} - \mathcal{H}_{\mathcal{C}-\{j\}}^{l,\varepsilon} * \overline{S}_{\mathcal{C}-\{j\}} \right)(t), \forall t : t \in [0, l] \right\}. \quad (121)$$

Eqn. (119) merely applies Definition 2, and Eqs. (120) and (121) use the definition of the convolution operator. For any fixed $t \geq 0$, setting $l = t$ in Eqs. (119)–(121) gives the claim

$$1 - \varepsilon \leq Pr \left\{ D_j(t) \geq A_j * \left(S_{\mathcal{C}} - \mathcal{H}_{\mathcal{C}-\{j\}}^{t,\varepsilon} * \overline{S}_{\mathcal{C}-\{j\}} \right)(t) \right\}. \quad (122)$$

Finally, since $D_j(t) \geq 0$ with probability one, the theorem follows. \square

D Proof of Theorem 4

We will verify that T_n defined by Eqn. (37) satisfies Eqn. (38).

Consider first the case $n = 1$. If $T_1 \geq T_o$, there is nothing to show since then Eqn. (36) implies Eqn. (38). If $T_1 < T_o$, then, by definition of T_1 , there exists a sequence τ_i in $[T_1, T_o]$ with $\lim_{i \rightarrow \infty} \tau_i = T_1$ so that

$$\mathcal{H}_{\mathcal{C}}^{2T_o, \varepsilon}(\tau_i) \leq S_{\mathcal{C}}(\tau_i) \quad \text{for all } i \geq 1. \quad (123)$$

Now, fix $t \geq 0$, and let $\underline{t}_{\mathcal{C}}$ and $\bar{t}_{\mathcal{C}}$ be the endpoints of the busy period as defined in Eqs. (33) and (34). By Eqs. (35) and (36), we have that $\underline{t}_{\mathcal{C}} \geq t - T_o$ and $\bar{t}_{\mathcal{C}} \leq t + T_o$.

By definition of $\mathcal{H}_{\mathcal{C}}^{2T_o, \varepsilon}$,

$$Pr \left\{ A_{\mathcal{C}}(\underline{t}_{\mathcal{C}}, \underline{t}_{\mathcal{C}} + \tau) \leq \mathcal{H}_{\mathcal{C}}^{2T_o, \varepsilon}(\tau), \forall \tau \leq T_o \right\} \geq 1 - \varepsilon. \quad (124)$$

Thus,

$$Pr \left\{ A_{\mathcal{C}}(\underline{t}_{\mathcal{C}}, \underline{t}_{\mathcal{C}} + \tau) - S_{\mathcal{C}}(\tau) \leq \mathcal{H}_{\mathcal{C}}^{2T_o, \varepsilon}(\tau) - S_{\mathcal{C}}(\tau), \forall \tau \leq T_o \right\} \geq 1 - \varepsilon. \quad (125)$$

Combining Eqn. (123) with Eqn. (125), we obtain

$$Pr \left\{ A_{\mathcal{C}}(\underline{t}_{\mathcal{C}}, \underline{t}_{\mathcal{C}} + \tau_i) - S_{\mathcal{C}}(\tau_i) \leq 0 \text{ for all } i \geq 1 \right\} \geq 1 - \varepsilon. \quad (126)$$

We conclude that with probability at least $1 - \varepsilon$, there must be times $x_i \leq \tau_i$, such that $A_{\mathcal{C}}(\underline{t}_{\mathcal{C}}, \underline{t}_{\mathcal{C}} + x_i) - D_{\mathcal{C}}(\underline{t}_{\mathcal{C}}, \underline{t}_{\mathcal{C}} + x_i) \leq 0$, that is, where the backlog is zero. Since $\lim_{i \rightarrow \infty} \tau_i = T_1$, we have with $1 - \varepsilon$ that $\bar{t}_{\mathcal{C}} - \underline{t}_{\mathcal{C}} \leq T_1$, which settles the case $n = 1$. The claim for $n > 1$ can be shown by induction. \square