

Extending the Network Calculus Pay Bursts Only Once Principle to Aggregate Scheduling

Markus Fidler

Department of Computer Science, Aachen University
Ahornstr. 55, 52074 Aachen, Germany
fidler@i4.informatik.rwth-aachen.de

Abstract. The Differentiated Services framework allows to provide scalable network Quality of Service by aggregate scheduling. Services, like a Premium class, can be defined to offer a bounded end-to-end delay. For such services, the methodology of Network Calculus has been applied successfully in Integrated Services networks to derive upper bounds on the delay of individual flows. Recent extensions allow an application of Network Calculus even to aggregate scheduling networks. Nevertheless, computations are significantly complicated due to the multiplexing and de-multiplexing of micro-flows to aggregates. Here problems concerning the tightness of delay bounds may be encountered.

A phenomenon called Pay Bursts Only Once is known to give a closer upper estimate on the delay, when an end-to-end service curve is derived prior to delay computations. Doing so accounts for bursts of the flow of interest only once end-to-end instead of at each link independently. This principle also holds in aggregate scheduling networks. However, it can be extended in that bursts of interfering flows are paid only once, too. In this paper we show the existence of such a complementing Pay Bursts Only Once phenomenon for interfering flows. We derive the end-to-end service curve for a flow of interest in an arbitrary aggregate scheduling feed forward network for rate-latency service curves, and leaky bucket constraint arrival curves, which conforms to both of the above principles. We give simulation results to show the utility of the derived forms.

1 Introduction

The Differentiated Services (DS) architecture [2] is the most recent approach of the Internet Engineering Task Force (IETF) towards network Quality of Service (QoS). DS addresses the scalability problems of the former Integrated Services approach by an aggregation of micro-flows to a small number of different traffic classes, for which service differentiation is provided. Packets are identified by simple markings that indicate the respective class. In the core of the network, routers do not need to determine to which flow a packet belongs, only which aggregate behavior has to be applied. Edge routers mark packets and indicate whether they are within profile or, if they are out of profile, in which case they might even be discarded by a dropper at the edge router. A particular marking on

a packet indicates a so-called Per Hop Behavior (PHB) that has to be applied for forwarding of the packet. Currently, the Expedited Forwarding (EF) PHB [9], and the Assured Forwarding (AF) PHB group are specified. The EF PHB is intended for building a service that offers low loss, low delay, and low delay jitter, namely a Premium service. The specification of the EF PHB was recently redefined to allow for a more exact and quantifiable definition [5]. Especially the derivation of delay bounds is of interest, when providing a Premium service. In [4] such bounds are derived for a general topology and a maximum load. However, these bounds can be improved, when additional information concerning the current load, and the special topology of a certain DS domain is available.

In [15] a central resource management for DS domains called a Bandwidth Broker (BB) is presented. A BB is a middleware service which controls and facilitates the dynamic access to network services of a particular administrative domain [10]. The task of a BB in a DS domain is to perform a careful admission control, and to set up the appropriate configuration of the domain's edge routers, whereas the configuration of core routers is intended to remain static to allow for scalability. While doing so, the BB knows about all requests for capacity of certain QoS classes. Besides it can easily learn about the DS domains topology, either statically, or by implementing a listener for the domains routing protocol. Thus, a BB can have access to all information that is required, to apply the mathematical methodology of Network Calculus [3, 13], in order to base its admission control on delay boundaries that are derived for the current load, and the special topology of the administrated domain [19].

In this paper we address the derivation of end-to-end delay guarantees based on Network Calculus. We derive a closed form solution for the end-to-end delay in feed forward First In First Out (FIFO) networks, for links that have a rate-latency property, and flows that are leaky bucket constraint. In particular this form accounts for bursts of interfering flows only once, and thus implements a principle for aggregate scheduling that is similar to the Network Calculus Pay Bursts Only Once principle [13]. The derived boundaries can be applied as a decision criterion to perform the admission control of a DS domain by a BB. The remainder of this paper is organized as follows: In Section 2 the required background on Network Calculus, and the notation that is applied in the sequel are given. Section 3 introduces two examples, which show how bursts of interfering flows worsen the derived delay bounds and, which prove the existence of a counterpart to the Pay Bursts Only Once phenomenon for interfering flows in aggregate scheduling networks. The first example can be satisfactorily solved by direct application of current Network Calculus, whereas to our knowledge for the second example a tight solution is missing in current literature. This missing piece is addressed in Section 4, where a tight closed form solution for arbitrary feed forward networks with FIFO rate-latency service curves and leaky bucket constraint arrival curves is derived. In Section 5 we describe the implementation of the admission control in a DS BB that is based on worst-case delay bounds. Numerical results on the performance gain that is achieved by applying the previously derived terms are given. Section 6 concludes the paper.

2 Network Calculus Background and Notation

Network Calculus is a theory of deterministic queuing systems that is based on the early work on the calculus for network delay in [6, 7], and on the work on Generalized Processor Sharing (GPS) presented in [16, 17]. Further extensions, and a comprehensive overview on current Network Calculus are given in [12, 13], and from the perspective of filtering theory in [3]. Here only a few concepts are covered briefly, to give the required background, and to introduce the notation that is mainly taken over from [13]. In addition since networks consisting of several links n that are used by a flow of interest, and a number of interfering flows m are investigated, upper indices j indicate links, and lower indices i indicate flows in the sequel.

The scheduler on an outgoing link can be characterized by the concept of a service curve, denoted by $\beta(t)$. A special characteristic of a service curve is the rate-latency type that is given by $\beta_{R,T}(t) = R \cdot [t - T]^+$ with a rate R and a latency T . The term $[x]^+$ is equal to x , if $x \geq 0$, and zero otherwise. Service curves of the rate-latency type are implemented for example by Priority Queuing (PQ), or Weighted Fair Queuing (WFQ). The latency of a PQ scheduler is given in [5] for variable length packet networks with a Maximum Transmission Unit (MTU) according to $T = \text{MTU}/R$. Nevertheless, routers can implement additional non-preemptive layer 2 queues on their outgoing interfaces for a smooth operation, which can add further delay to a layer 3 QoS implementation [20]. Thus $T = (l_2 + 1) \cdot \text{MTU}/R$ might have to be applied, whereby l_2 gives the layer 2 queuing capacity in units of the MTU.

Flows are defined either by their arrival functions denoted by $F(t)$, or by their arrival curves $\alpha(t)$, whereas $\alpha(t_2 - t_1) \geq F(t_2) - F(t_1)$ for all $t_2 \geq t_1$. In DS networks, a typical characteristic for incoming flows can be given by the leaky bucket constraint $\alpha_{r,b}(t) = b + r \cdot t$ that is also known as sigma-rho leaky bucket in [3]. Usually the ingress router of a DS domain meters incoming flows against a leaky bucket algorithm, and either shapes, or drops non-conforming traffic, which justifies the application of leaky bucket constraint arrival curves.

If a link j is traversed by a flow i , the arrival function of the output flow F_i^{j+1} , which is the input arrival function for an existing, or an imaginary subsequent link $j + 1$, can be given according to (1) for $t \geq s \geq 0$ [12].

$$F_i^{j+1}(t) \geq F_i^j(t - s) + \beta^j(s) \quad (1)$$

From (1) the term in (2) follows. The operator \otimes denotes the convolution under the min-plus algebra that is applied by Network Calculus.

$$F_i^{j+1}(t) \geq (F_i^j \otimes \beta^j)(t) = \inf_{t \geq s \geq 0} [F_i^j(t - s) + \beta^j(s)] \quad (2)$$

Further on, the output flow is upper constrained by an arrival curve α_i^{j+1} that is given according to (3), with \oslash denoting the min-plus de-convolution.

$$\alpha_i^{j+1}(t) = (\alpha_i^j \oslash \beta^j)(t) = \sup_{s \geq 0} [\alpha_i^j(t + s) - \beta^j(s)] \quad (3)$$

If the path of a flow i consists of two or more links, the formulation of the concatenation of links can be derived based on (4).

$$\frac{F_i^{j+2}(u) - F_i^{j+1}(u - (t - s))}{F_i^{j+2}(u)} \geq \frac{\beta^{j+1}(t - s)}{\beta^j(s)} \quad (4)$$

The end-to-end service curve is then given in (5), which covers the case of two links, whereas the direct application of (5) also holds for n links.

$$\beta^{j+1,j}(t) = (\beta^{j+1} \otimes \beta^j)(t) = \inf_{t \geq s \geq 0} [\beta^{j+1}(t - s) + \beta^j(s)] \quad (5)$$

The maximal virtual delay d for a system that offers a service curve of $\beta(t)$ with an input flow that is constraint by $\alpha(t)$, is given as the supremum of the horizontal deviation according to (6).

$$d \leq \sup_{s \geq 0} [\inf_{\tau \geq 0 : \alpha(s) \leq \beta(s + \tau) + \beta^j(s)} \tau] \quad (6)$$

For the special case of service curves of the rate-latency type, and sigma-rho leaky bucket constraint arrival curves, simplified solutions exist for the above equations. The arrival curve of the output flow according to (3) is given in (7) for this case, whereas the burst size of the output flow $b_i + r_i \cdot T^j$ is equal to the maximum backlog at the scheduler of the outgoing link.

$$\alpha_i^{j+1}(t) = b_i + r_i \cdot T^j + r_i \cdot t \quad (7)$$

The concatenation of two rate-latency service curves can be reduced to (8).

$$\beta^{j+1,j}(t) = \min[R^{j+1}, R^j] \cdot [t - (T^{j+1} + T^j)]^+ \quad (8)$$

Finally, (9) gives the worst case delay for the combination of a rate-latency service curve, and a leaky bucket constraint arrival curve.

$$d \leq T^j + b_i/R^j \quad (9)$$

If a flow i traverses two links j and $j + 1$, two options for the derivation of the end-to-end delay exist. For simplicity, the service curves are assumed here to be of the rate-latency type, and the arrival curve is chosen to be leaky bucket constraint. At first the input arrival curve of flow i at link $j + 1$ can be computed as in (7). Then the virtual end-to-end delay is derived to be the sum of the virtual delays at link j and $j + 1$ according to (9). Doing so results in $d \leq T^j + b_i/R^j + T^{j+1} + (b_i + r_i \cdot T^j)/R^{j+1}$. The second option is to derive the end-to-end service curve as is done in (8), and compute the delay according to (9) afterwards, resulting in $d \leq T^j + T^{j+1} + b_i/\min[R^{j+1}, R^j]$. Obviously, the second form gives a closer bound, since it accounts for the burst size b_i only once. This property is known as the Pay Bursts Only Once phenomenon from [13].

Until now only networks that perform a per-flow based scheduling, for example Integrated Services networks, have been considered. The aggregation or

the multiplexing of flows can be given by the addition of the arrival functions $F_{1,2}(t) = F_1(t) + F_2(t)$, or arrival curves $\alpha_{1,2}(t) = \alpha_1(t) + \alpha_2(t)$. For aggregate scheduling networks with FIFO service curves, families of per-flow service curves $\beta_\theta(t)$ according to (10) with an arbitrary parameter $\theta \geq 0$ are derived in [8, 13]. $\beta_\theta^j(t)$ gives a family of service curves for a flow 1 that is scheduled in an aggregate manner in conjunction with a flow 2 on a link j . $1_{t>\theta}$ is zero for $t \leq \theta$.

$$\beta_\theta^j(t) = [\beta^j(t) - \alpha_2(t - \theta)]^+ 1_{t>\theta} \quad (10)$$

The parameter θ has to be set to zero in case of blind multiplexing. In FIFO networks for $u = \sup[v : F_{1,2}^j(v) \leq F_{1,2}^{j+1}(t)]$, the conditions $F^j(u) \leq F^{j+1}(t)$, and $F^j(u^+) \geq F^{j+1}(t)$ with $u^+ = u + \epsilon$, $\epsilon > 0$ hold for the sum of both flows, and for the individual flows, too. These additional constraints allow to derive (10) for an arbitrary $\theta > 0$. The complete derivation is missed out here. It can be found in [13]. From (10) it cannot be concluded that $\sup_\theta[\beta_\theta^j]$, or $\inf_\theta[\beta_\theta^j]$ is a service curve, but for the output flow $\alpha_1^{j+1}(t) = \inf_{\theta \geq 0}[(\alpha_1^j \circ \beta_\theta^j)(t)]$ can be given [13].

For the special case of rate-latency service curves, and leaky bucket constraint arrival curves, (11) can be derived from (10) to be a service curve for flow 1 for $r_1 + r_2 < R^j$.

$$\beta^j(t) = (R^j - r_2) \cdot [t - (T^j + b_2/R^j)]^+ \quad (11)$$

The methodology that is shown up to here already allows to implement the algorithmic derivation of delay bounds in a DS domain by a BB, with the restriction that the domain has to be a feed forward network. In non feed forward networks analytical methods like time stopping [3, 13] can be applied. Nevertheless, for an algorithmic implementation, it is simpler to prevent from aggregate cycles, and transform the domains topology into a feed forward network, to allow for the direct application of Network Calculus. Such a conversion of the network topology can be made by means of breaking up loops by forbidding certain links, for example by turn prohibition as presented in [21] for networks that consist of bidirectional links. Doing so, the BB can inductively derive the arrival curves of micro-flows at each outgoing link, then compute the service curve of each link from the point of view of each micro-flow, and concatenate these to end-to-end service curves, to give upper bounds on the worst case delay for individual flows.

3 Extended Pay Bursts Only Once Principle

Though Section 2 gives the required background to derive per micro-flow based end-to-end delay bounds in aggregate scheduling networks, these bounds are likely to be unnecessarily loose. Figure 1 gives a motivating example of a simple network consisting of two outgoing links that are traversed by two flows. This example is given to introduce a similar concept to the Pay Bursts Only Once principle [13] for aggregate scheduling. Flow 1 is the flow of interest, for which the end-to-end delay needs to be derived. According to the Pay Bursts Only Once principle, at first the end-to-end service curve for flow 1 has to be derived, and then the delay is computed, instead of computing the delay at each link

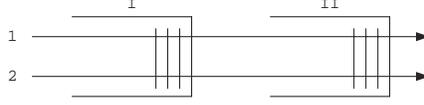


Fig. 1. Two Flows Share Two Consecutive Links

independently, and summing these delays up. A similar decision has to be made, when deriving the service curve for flow 1 by subtracting the arrival curve of flow 2 from the aggregate service curves. Either this subtraction is made at each link independently before concatenating service curves, or the subtraction is made after service curves have been concatenated. To illustrate the difference the two possible derivations of the end-to-end service curve for flow 1 are given here for rate-latency service curves, and leaky bucket constraint arrival curves.

For the first option, which is already drafted at the end of Section 2, the service curve of link I from the point of view of flow 2 can be derived to be of the rate-latency type with a rate $R^I - r_1$ and a latency $T^I + b_1/R^I$. Then, the arrival curve of flow 2 at link II can be given to be leaky bucket constraint with the rate r_2 and the burst size $b_2 + r_2 \cdot (T^I + b_1/R^I)$. The service curves at link I and II from the point of view of flow 1 can be given to be rate-latency service curves with the rates $R^I - r_2$, respective $R^{II} - r_2$, and the latencies $T^I + b_2/R^I$, respective $T^{II} + (b_2 + r_2 \cdot (T^I + b_1/R^I))/R^{II}$. The concatenation of these two service curves yields the rate-latency service curve for flow 1 given in (12).

$$\beta^{I,II}(t) = \min[R^I - r_2, R^{II} - r_2] \cdot [t - (T^I + b_2/R^I) - (T^{II} + (b_2 + r_2 \cdot (T^I + b_1/R^I))/R^{II})]^+ \quad (12)$$

The second option requires that the service curves of link I and II are convoluted prior to subtraction of the flow 2 arrival curve [19]. The concatenation yields a service curve of the rate-latency type with a rate of $\min[R^I, R^{II}]$ and a latency of $T^I + T^{II}$. After subtracting the arrival curve of flow 2 from the concatenation of the service curves of link I, and II, (13) can be given for flow 1.

$$\beta^{I,II}(t) = (\min[R^I, R^{II}] - r_2) \cdot [t - (T^I + T^{II} + b_2/\min[R^I, R^{II}])]^+ \quad (13)$$

Obviously, the form in (13) offers a lower latency than the one in (12), which accounts for the burst size of the interfering flow 2 twice, once with the size b_2 at link I, and then with $b_2 + r_2 \cdot (T^I + b_1/R^I)$ at link II. The increase of the burst size of flow 2 at link II is due to the aggregate scheduling of flow 2 with flow 1 at link I. Thus, based on (12), and (13) a counterpart to the Pay Burst Only Once phenomenon has been shown for interfering flows in aggregate scheduling networks.

However, most problems of interest cannot be solved as simple as the one in Figure 1. One such example is shown in Figure 2. Flow 2 is the flow of interest, for which the end-to-end service curve needs to be derived. Links I and II can

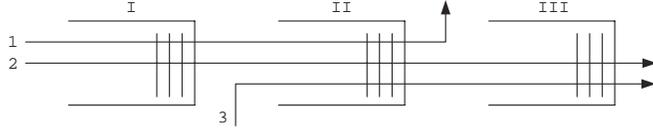


Fig. 2. Three Flows Share Three Consecutive Links

be concatenated after the arrival curve of flow 3 is subtracted from the service curve of link II. Then the arrival curve of flow 1 can be subtracted from the concatenation of the service curves of link I and II. Unfortunately the direct application of the Network Calculus terms given in Section 2 then requires that the arrival curve of flow 3 at link III is subtracted from the service curve of link III independently, which violates the extended Pay Bursts Only Once principle. An alternative derivation is possible, if the arrival curve of flow 1 is subtracted from the service curve of link II before links II and III are concatenated. Doing so does unfortunately encounter the same problem.

4 Closed Form Solution for Feed Forward Networks

Consider an end-to-end path of a flow of interest i in an arbitrary aggregate scheduling feed forward network with FIFO service curve elements. Further on, assume that a concatenation of consecutive links offers a FIFO characteristic, too. The path consists of n links that are indexed by j in ascending order from the source to the destination of the flow of interest i . These links are given in the set \mathbb{J}_i . Further on, the links of the path of flow i are used in an aggregate manner by an additional number of flows m that are indexed by k . The paths of the flows k are given by the sets \mathbb{J}_k . For each link j a set \mathbb{K}^j is defined to hold all other flows k that traverse the link, not including flow i . Note that flows may exist that share a part of the path of the flow of interest i , then follow a different path, and afterwards again share a part of the path of flow i . These flows are split up in advance into as many individual flows as such multiplexing points exist. Thus, all resulting flows do have only one multiplexing point with flow i . In [13] the term Route Interference Number (RIN) is defined to hold the quantity of such multiplexing points. Then, the set $\mathbb{K}_i = \bigcup_{j \in \mathbb{J}_i} \mathbb{K}^j$ is defined to hold all m flows that use the complete path or some part of the path of flow i . Further on, the set $\mathbb{J}_{i,k} = \mathbb{J}_i \cap \mathbb{J}_k$ holds all links of a sub-path that are used by both flow i and a flow k . Based on these definitions, on the terms in (4), and on the rules for multiplexing, (14) can be given for time pairs $t_{j+1} - t_j \geq 0$ for all $j \in \mathbb{J}_i$. The time indices are chosen here to match link indices, since arrival functions are observed at these time instances at the belonging links.

$$F_i^{n+1}(t_{n+1}) - F_i^1(t_1) \geq \sum_{j \in \mathbb{J}_i} \beta^j(t_{j+1} - t_j) - \sum_{k \in \mathbb{K}_i} \sum_{j \in \mathbb{J}_{i,k}} (F_k^{j+1}(t_{j+1}) - F_k^j(t_j)) \quad (14)$$

Now, the arrival functions of the interfering flows can be easily split up and subtracted from the relevant links. Like for the derivation of (10) in [13], FIFO conditions can be defined on a per-link basis. Doing so, pairs of arrival functions $F_k^{j+1}(t_{j+1}) - F_k^j(t_j)$ can be replaced by their arrival curves $\alpha_k^j(t_{j+1} - t_j - \theta^j)$ for $t_{j+1} - t_j > \theta^j$, with the arbitrary per link parameters $\theta^j \geq 0$. Nevertheless, doing so results in paying the bursts of all interfering flows at each link independently.

However, FIFO conditions can in addition to per link be derived per sub-path $\mathbb{J}_{i,k}$. With $\sum_{j \in \mathbb{J}_{i,k}} (F_k^{j+1}(t_{j+1}) - F_k^j(t_j)) = F_k^{j_{\max}+1}(t_{j_{\max}+1}) - F_k^{j_{\min}}(t_{j_{\min}})$, whereby $j_{\max} = \max[j \in \mathbb{J}_{i,k}]$, and $j_{\min} = \min[j \in \mathbb{J}_{i,k}]$, (15) can be derived for $t_{j_{\max}+1} - t_{j_{\min}} > \theta_k$, based on (10).

$$F_i^{n+1}(t_{n+1}) - F_i^1(t_1) \geq \sum_{j \in \mathbb{J}_i} \beta^j(t_{j+1} - t_j) - \sum_{k \in \mathbb{K}_i} \alpha_k^{j_{\min}}(t_{j_{\max}+1} - t_{j_{\min}} - \theta_k) \quad (15)$$

To motivate the step from (14) to (15), the FIFO conditions that are applied, are shown exemplarily for the small network in Figure 2. The derivation, is a direct application of the one given in [13] for the form in (10). Define an u_1 , and u_3 for flow 1, respective flow 3 according to (16), and (17).

$$u_1 = \sup\{v : F_1^I(v) + F_2^I(v) \leq F_1^{III}(t_3) + F_2^{III}(t_3)\} \quad (16)$$

$$u_3 = \sup\{v : F_2^{II}(v) + F_3^{II}(v) \leq F_2^{IV}(t_4) + F_3^{IV}(t_4)\} \quad (17)$$

Now, from (16), and with $u_1^+ = \inf\{v : v > u_1\}$, $t_1 \leq u_1 \leq t_3$, $F_1^I(u_1) + F_2^I(u_1) \leq F_1^{III}(t_3) + F_2^{III}(t_3)$, and $F_1^I(u_1^+) + F_2^I(u_1^+) \geq F_1^{III}(t_3) + F_2^{III}(t_3)$ can be given. Due to the per sub-path FIFO conditions, the above terms also hold for the individual flows 1, and 2, for example $F_1^I(u_1) \leq F_1^{III}(t_3)$, and $F_1^I(u_1^+) \geq F_1^{III}(t_3)$. Similar conditions can be derived for flows 2, and 3 from (17). The following term in (18) can be set up based on (4).

$$F_2^{IV}(t_4) - F_2^I(t_1) \geq \beta^{III}(t_4 - t_3) + \beta^{II}(t_3 - t_2) + \beta^I(t_2 - t_1) \\ - (F_3^{IV}(t_4) - F_3^{II}(t_2)) - (F_1^{III}(t_3) - F_1^I(t_1)) \quad (18)$$

With $F_1^I(u_1^+) \geq F_1^{III}(t_3)$, and $F_3^{II}(u_3^+) \geq F_3^{IV}(t_4)$, (19) can be derived.

$$F_2^{IV}(t_4) - F_2^I(t_1) \geq \beta^{III}(t_4 - t_3) + \beta^{II}(t_3 - t_2) + \beta^I(t_2 - t_1) \\ - (F_3^{II}(u_3^+) - F_3^{II}(t_2)) - (F_1^I(u_1^+) - F_1^I(t_1)) \quad (19)$$

Choose two arbitrary parameters $\theta_1 \geq 0$, and $\theta_3 \geq 0$. If $u_1 < t_3 - \theta_1$, $F_1^I(u_1^+) \leq F_1^I(t_3 - \theta_1)$ holds. Further on, $t_3 - t_1 > \theta_1$ can be given in this case, since $u_1 \geq t_1$. With a similar condition for flow 3, the form in (20) is derived.

$$F_2^{IV}(t_4) - F_2^I(t_1) \geq \beta^{III}(t_4 - t_3) + \beta^{II}(t_3 - t_2) + \beta^I(t_2 - t_1) \\ - \alpha_3(t_4 - t_2 - \theta_3) - \alpha_1(t_3 - t_1 - \theta_1) \quad (20)$$

Else, for $u_1 \geq t_3 - \theta_1$, the FIFO condition $F_2^{III}(t_3) \geq F_2^I(u_1)$ is applied. Further on, the term $F_2^{III}(t_3) - F_2^I(t_1) \geq \beta^{I,II}(t_3 - t_1)$ can be set up, with $\beta^{I,II}$ denoting

the service curve for flow 2 that is offered by link I, and II. Substitution of t_1 by u_1 , yields $F_2^{\text{III}}(t_3) \geq F_2^{\text{I}}(u_1) + \beta^{\text{I,II}}(t_3 - u_1)$. Hence, $\beta^{\text{I,II}}(t_3 - t_1) = 0$ is a trivial service curve for $u_1 \geq t_3 - \theta_1$. Similar forms can be derived for $u_3 \geq t_4 - \theta_3$.

In the following the form in (15) is solved for the simple case of rate-latency service curves, and leaky bucket constraint arrival curves.

Proposition 1 (End-to-End Service Curve) *The end-to-end service curve for a flow of interest i in an aggregate scheduling feed forward network with FIFO service curve elements of the rate-latency type $\beta_{R,T}$, and leaky bucket constrained arrival curves $\alpha_{r,b}$ is again of the rate-latency type, and given according to (21).*

$$\beta_i(t) = \min_{j \in \mathbb{J}_i} \left[R^j - \sum_{k \in \mathbb{K}^j} r_k \right] \cdot \left[t - \sum_{j \in \mathbb{J}_i} T^j - \sum_{k \in \mathbb{K}_i} \frac{b_k^{j_{\min}}}{\min_{j \in \mathbb{J}_{i,k}} [R^j]} \right]^+ \quad (21)$$

The form in (21) gives an intuitive result. The end-to-end service curve for flow i has a rate R , which is the minimum of the remaining rates at the traversed links, after subtracting the rates of interfering flows that share the individual links. The latency T is given as the sum of all latencies along the path, plus the burst size of interfering flows at their multiplexing points indicated by j_{\min} , divided by the minimum rate along the common sub-path with the flow of interest i . Thus, bursts of interfering flows account only once with the maximal latency that can be evoked by such bursts along the shared sub-paths.

Proof 1 In (22) the form in (15) is given for rate-latency service curves, and leaky bucket constraint arrival curves for $t_{j_{\max}+1} - t_{j_{\min}} > \theta_k$.

$$\begin{aligned} F_i^{n+1}(t_{n+1}) - F_i^1(t_1) &\geq \sum_{j \in \mathbb{J}_i} (R^j \cdot [t_{j+1} - t_j - T^j]^+) \\ &\quad - \sum_{k \in \mathbb{K}_i} (b_k^{j_{\min}} + r_k \cdot (t_{j_{\max}+1} - t_{j_{\min}} - \theta_k)) \end{aligned} \quad (22)$$

Here, we apply a definition of the per flow θ_k by per link θ^j according to (23). Now the failure of any of the conditions $t_{j_{\max}+1} - t_{j_{\min}} > \theta_k$ requires the failure of at least one of the conditions $t_{j+1} - t_j > \theta^j$ for any j with $j_{\max} \geq j \geq j_{\min}$.

$$\theta_k = \sum_{j \in \mathbb{J}_{i,k}} \theta^j \quad (23)$$

Reformulation of (22) yields (24) for $t_{j+1} - t_j > \theta_j$, with $j \in \mathbb{J}_{i,k}$.

$$\begin{aligned} F_i^{n+1}(t_{n+1}) - F_i^1(t_1) &\geq \sum_{j \in \mathbb{J}_i} (R^j \cdot [t_{j+1} - t_j - T^j]^+) \\ &\quad - \sum_{k \in \mathbb{K}_i} \left(b_k^{j_{\min}} + r_k \cdot \sum_{j \in \mathbb{J}_{i,k}} (t_{j+1} - t_j - \theta^j) \right) \end{aligned} \quad (24)$$

Next, the fact that (10) gives a service curve for any setting of the parameter θ with $\theta \geq 0$ is used. Thus, the per-flow θ_k can be set arbitrarily with $\theta_k \geq$

0, whereas this condition can according to (23) be fulfilled by any $\theta^j \geq 0$. Hence, (24) gives a service curve for any $\theta^j \geq 0$. We define an initial setting of the parameters θ^j , for which an end-to-end service curve for flow i is derived for all $t_{j+1} - t_j > \theta^j$. Then, for the special cases in which $t_{j+1} - t_j > \theta^j$ does not hold for one or several links j , θ^j is redefined. We show for these cases by means of the redefined θ^j that the end-to-end service curve that is derived before still holds true. Thus, we prove that this service curve is valid for all $t_{j+1} - t_j \geq 0$. The definition of different settings of the parameters θ^j is not generally allowed to derive a service curve. As already stated in Section 2, it cannot be concluded that for example $\inf_{\theta} \beta_{\theta}(t)$ is a service curve [13]. Nevertheless, there is a difference between doing so, and the derivation that is shown in the following. Here, a service curve β_{θ^j} is derived for fixed θ^j for all $t_{j+1} - t_j > \theta^j$. This service curve is not modified later on, but only proven to hold for any $t_{j+1} - t_j \geq 0$ by applying different settings of the θ^j . The initially applied setting of the θ^j is given in (25). The θ^j are defined to be the latency of the scheduler on the outgoing link j plus the burst size of interfering flows k , if j is the multiplexing point of the flow k with flow i , divided by the minimum rate along the common sub-path of flow k and i .

$$\theta^j = T^j + \sum_{k \in \mathbb{K}^j | j = j_{\min}} \frac{b_k^{j_{\min}}}{\min_{j' \in \mathbb{J}_{i,k}} [R^{j'}]} \quad (25)$$

With (25) the term in (24) can be rewritten according to (26) for all $t_{j+1} - t_j > \theta^j$.

$$\begin{aligned} F_i^{n+1}(t_{n+1}) - F_i^1(t_1) &\geq \sum_{j \in \mathbb{J}_i} \left(R^j \cdot [t_{j+1} - t_j - T^j]^+ \right) \\ &- \sum_{k \in \mathbb{K}_i} \left(b_k^{j_{\min}} + r_k \cdot \sum_{j \in \mathbb{J}_{i,k}} \left(t_{j+1} - t_j - T^j - \sum_{k' \in \mathbb{K}^j | j = j_{\min}} \frac{b_{k'}^{j_{\min}}}{\min_{j' \in \mathbb{J}_{i,k'}} [R^{j'}]} \right) \right) \end{aligned} \quad (26)$$

Some reordering, while scaling up the subtrahends by adding $[\dots]^+$ conditions, and a replacement of $\sum_{k \in \mathbb{K}_i} \sum_{j \in \mathbb{J}_{i,k}}$ by $\sum_{j \in \mathbb{J}_i} \sum_{k \in \mathbb{K}^j}$ yields (27).

$$\begin{aligned} F_i^{n+1}(t_{n+1}) - F_i^1(t_1) &\geq \sum_{j \in \mathbb{J}_i} \left((R^j - \sum_{k \in \mathbb{K}^j} r_k) \cdot [t_{j+1} - t_j - T^j]^+ \right) \\ &- \sum_{k \in \mathbb{K}_i} \left(b_k^{j_{\min}} - r_k \cdot \sum_{j \in \mathbb{J}_{i,k}} \sum_{k' \in \mathbb{K}^j | j = j_{\min}} \frac{b_{k'}^{j_{\min}}}{\min_{j' \in \mathbb{J}_{i,k'}} [R^{j'}]} \right) \end{aligned} \quad (27)$$

With the replacement of $\sum_{k \in \mathbb{K}_i} \sum_{j \in \mathbb{J}_{i,k}}$ by $\sum_{j \in \mathbb{J}_i} \sum_{k \in \mathbb{K}^j}$, (28) can be derived.

$$\begin{aligned} F_i^{n+1}(t_{n+1}) - F_i^1(t_1) &\geq \sum_{j \in \mathbb{J}_i} \left((R^j - \sum_{k \in \mathbb{K}^j} r_k) \cdot [t_{j+1} - t_j - T^j]^+ \right) \\ &- \sum_{k \in \mathbb{K}_i} b_k^{j_{\min}} + \sum_{j \in \mathbb{J}_i} \sum_{k \in \mathbb{K}^j} \left(r_k \cdot \sum_{k' \in \mathbb{K}^j | j = j_{\min}} \frac{b_{k'}^{j_{\min}}}{\min_{j' \in \mathbb{J}_{i,k'}} [R^{j'}]} \right) \end{aligned} \quad (28)$$

Further on $\sum_{k \in \mathbb{K}_i} b_k^{j \min} = \sum_{j \in \mathbb{J}_i} \sum_{k' \in \mathbb{K}^j | j = j_{\min}} b_{k'}^{j \min}$ yields (29).

$$F_i^{n+1}(t_{n+1}) - F_i^1(t_1) \geq \sum_{j \in \mathbb{J}_i} \left((R^j - \sum_{k \in \mathbb{K}^j} r_k) \cdot [t_{j+1} - t_j - T^j]^+ \right) - \sum_{j \in \mathbb{J}_i} \left(\sum_{k' \in \mathbb{K}^j | j = j_{\min}} b_{k'}^{j \min} - \sum_{k \in \mathbb{K}^j} r_k \cdot \sum_{k' \in \mathbb{K}^j | j = j_{\min}} \frac{b_{k'}^{j \min}}{\min_{j' \in \mathbb{J}_{i,k'}} [R^{j'}]} \right) \quad (29)$$

Applying the common denominator, while scaling up the subtrahend, and with $j \in \mathbb{J}_{i,k'}$ and thereby $R^j \geq \min_{j' \in \mathbb{J}_{i,k'}} [R^{j'}]$ with $k' \in \mathbb{K}^j$, (30) can be derived.

$$F_i^{n+1}(t_{n+1}) - F_i^1(t_1) \geq \sum_{j \in \mathbb{J}_i} \left((R^j - \sum_{k \in \mathbb{K}^j} r_k) \cdot [t_{j+1} - t_j - T^j]^+ \right) - \sum_{j \in \mathbb{J}_i} \left(\left(R^j - \sum_{k \in \mathbb{K}^j} r_k \right) \cdot \sum_{k' \in \mathbb{K}^j | j = j_{\min}} \frac{b_{k'}^{j \min}}{\min_{j' \in \mathbb{J}_{i,k'}} [R^{j'}]} \right) \quad (30)$$

Then, (30) can be reformulated according to (31), still for $t_{j+1} - t_j > \theta_j$.

$$F_i^{n+1}(t_{n+1}) - F_i^1(t_1) \geq \sum_{j \in \mathbb{J}_i} \left((R^j - \sum_{k \in \mathbb{K}^j} r_k) \cdot \left([t_{j+1} - t_j - T^j]^+ - \sum_{k' \in \mathbb{K}^j | j = j_{\min}} \frac{b_{k'}^{j \min}}{\min_{j' \in \mathbb{J}_{i,k'}} [R^{j'}]} \right) \right) \quad (31)$$

The $\inf_{(t_{j+1} - t_j > \theta_j) | j \in \mathbb{J}_i}$ of (31) can be derived to be the form that is given in (21). Thus, the service curve in (21) is approved for all $t_{j+1} - t_j > \theta^j$ with the parameter settings of θ^j according to (25).

Now, if some of the conditions $t_{j+1} - t_j > \theta^j$ fail, the θ^j that are given in (25) can be redefined according to (32), based on the arbitrary parameters $\delta_k^j \geq 0$ with $\sum_{j \in \mathbb{J}_{i,k}} \delta_k^j = 1$.

$$\theta^j = T^j + \sum_{k \in \mathbb{K}^j} \frac{\delta_k^j \cdot b_k^{j \min}}{\min_{j' \in \mathbb{J}_{i,k}} [R^{j'}]} \quad (32)$$

The burst size of interfering flows $b_k^{j \min}$ is arbitrarily accounted for by θ^j in (25), whereas, if an interfering flow k traverses more than one link, the burst size $b_k^{j \min}$ could be part of any θ^j , with $j \in \mathbb{J}_{i,k}$. For such redefined θ^j , it can be shown that the same derivation as above holds, resulting in (33). Again, applying the $\inf_{(t_{j+1} - t_j > \theta_j) | j \in \mathbb{J}_i}$ leads to the same form (21), as shown for (31) before.

$$F_i^{n+1}(t_{n+1}) - F_i^1(t_1) \geq \sum_{j \in \mathbb{J}_i} \left((R^j - \sum_{k \in \mathbb{K}^j} r_k) \cdot \left([t_{j+1} - t_j - T^j]^+ - \sum_{k' \in \mathbb{K}^j} \frac{\delta_{k'}^j \cdot b_{k'}^{j \min}}{\min_{j' \in \mathbb{J}_{i,k'}} [R^{j'}]} \right) \right) \quad (33)$$

However, there can be pairs of $t_{j+1} - t_j \geq 0$, for which no setting of the parameters δ_k^j according to (32) allows a redefinition of θ^j , for which $t_{j+1} - t_j > \theta^j$ for all $j \in \mathbb{J}_i$ can be achieved. A condition $t_{j+1} - t_j > \theta_j$ can fail for $\delta_k^j = 0$ for all $k \in \mathbb{K}_j$, if $t_{j+1} - t_j \leq T^j$, without violating any of the per-flow conditions $t_{j_{\max}+1} - t_{j_{\min}} > \theta_k$. In this case the terms that are related to link j in (33) are nullified immediately by the $[\dots]^+$ condition. Nevertheless, if some of the per flow conditions $t_{j_{\max}+1} - t_{j_{\min}} > \theta_k$ are violated for a number of flows $k \in \mathbb{L}_i$, the service curves of the sub-paths $\bigcup_{k \in \mathbb{L}_i} \mathbb{J}_{i,k}$ have according to the derivation of (10) to be set to zero. However, setting the service curves of sub-paths to zero is the same as setting the service curves of all links along these paths to zero. Regarding (33), it can be seen that any links for which the service curve is set to zero possibly increase the resulting rate of the service curve, compared to (21), whereas the resulting maximum latency is not influenced. This holds true for any θ^j according to (32), respective for any δ_k^j with $\sum_{j \in \mathbb{J}_{i,k}} \delta_k^j = 1$. Thus, also the case of $\delta_k^j = 0$ for all $k \in \mathbb{K}_i \setminus \{\mathbb{L}_i\}$, and $j \in \bigcup_{k \in \mathbb{L}_i} \mathbb{J}_{i,k}$ is covered. The latter setting ensures that bursts of flows that share part of the sub-paths that are set to zero, but that also traverse further links, are accounted for at these links. Then by scaling down the term $\min_{j \in \mathbb{J}_i \setminus \{\mathbb{J}_{i,k} | k \in \mathbb{L}_i\}} [R^j - \sum_{k \in \mathbb{K}^j} r_k]$ to $\min_{j \in \mathbb{J}_i} [R^j - \sum_{k \in \mathbb{K}^j} r_k]$, (21) is also a service curve for cases in which $t_{j+1} - t_j \leq \theta^j$, and finally holds for any $t_{j+1} - t_j \geq 0$ with $j \in \mathbb{J}_i$. \square

Finally, (34) gives a tight end-to-end service curve for flow 2 in Figure 2. As intended, the bursts of the interfering flows 1, and 3 are accounted for only once, with their initial burst size at the multiplexing, or route interference point.

$$\begin{aligned} \beta_2(t) = & \min[R^I - r_1, R^{II} - r_1 - r_3, R^{III} - r_3] \\ & \cdot \left[t - T^I - T^{II} - T^{III} - \frac{b_1}{\min[R^I, R^{II}]} - \frac{b_3}{\min[R^{II}, R^{III}]} \right]^+ \end{aligned} \quad (34)$$

5 Numerical Results

In this section we give numerical results on the derivation of edge-to-edge delay bounds in a DS domain, which can be efficiently applied for the definition of so-called Per Domain Behaviors (PDBs) [14]. We compare the options of applying either the Extended Pay Burst Only Once principle, the Pay Bursts Only Once principle, or none of the two principles.

We implemented an admission control for an application as a BB in a DS domain. The BB currently knows about the topology of its domain statically, whereas a routing protocol listener can be added. Requests for Premium capacity are sent via socket communication to the BB. The requests consist of a start, and an end time to allow for both immediate, and advance reservation, a Committed Information Rate (CIR), a Committed Burst Size (CBS), and a target maximum delay. Whenever the BB receives a new request, it computes the edge-to-edge delay for all requests that are active during the period of time of the new request, as described in Section 2. If none of the target maximum per-flow delays is violated, the new request is accepted, which otherwise is rejected.

For performance evaluation we implemented a simulator that generates such Premium resource requests. Sources and sinks are chosen uniformly from a predefined set. Start, and end times are modelled as negative exponentially distributed with a mean λ , respective μ , that is a mean of $\rho = \mu/\lambda$ requests are active concurrently. This modelling has been found to be appropriate for user sessions, for example File Transfer Protocol (FTP) sessions in [18]. The target delay, CIR, and CBS are used as uniformly distributed parameters for the simulations.

The topology that is used is shown in Figure 3. It consists of the level one, and level two nodes of the German Research Network (DFN) [1]. The level one nodes are core nodes. End systems are connected to the level two nodes that are edge nodes. In detail, we connect up to five sources and sinks to each of the level two nodes. Links are either Synchronous Transfer Mode (STM) 4, STM 16, or STM 64 connections. The link transmission delay is assumed to be 2 ms. Shortest Path First (SPF) routing is applied to minimize the number of hops along the paths. Further on, Turn Prohibition (TP) [21] is used, to ensure the required feed forward property of the network. Figure 3 shows how loops are broken within the level one mesh of the DFN topology by the TP algorithm. The nodes have been processed by TP in the order of their numbering. For example the turn (8; 7; 9) that is the turn from node 8 via node 7 to node 9 is prohibited, whereas for instance the turn (8; 7; 6), or simply the use of the link (8; 7) is permitted. For the DFN topology the SPF TP algorithm does only increase the length of one of the paths by one hop compared to SPF routing. Further on, the TP algorithm can be configured to prohibit turns that include links with a comparably low capacity with priority [21], as is shown in Figure 3, and applied by our BB. The Premium service is implemented based on PQ. Thus, service curves are of the rate-latency type with a latency set to the time it takes to transmit 3 MTU of 9.6 kB, to account for non-preemptive scheduling, due to packetization, and a router internal buffer for 2 Jumbo frames [20].

The performance measure that we apply is the ratio of accepted requests divided by the overall number of requests, and as an alternative the distribution of the derived delay bounds. Simulations have been run, until the 0.95 confidence interval of the acceptance ratio was smaller than 0.01. Initial-data deletion [11] has been applied to capture only the steady-state behavior, and the replication and deletion approach for means, that is for example shown in [11], was used. The results are given for different settings of the requested maximum delay, CBS, and CIR, and for a varying load $\rho = \mu/\lambda$ in Figure 4, and 5. In addition Figure 6 shows the fraction of flows for which a delay bound that is smaller than the delay given on the abscissa was derived.

As one of our main results we find a significant performance gain, when accounting for the Extended Pay Bursts Only Once phenomenon. The Pay Bursts Only Once principle alone allows to derive noticeable tighter delay bounds, based on edge-to-edge service curves, compared to an incremental delay computation, as described already in Section 2. The advantage can further on be seen in terms of the acceptance ratio in Figure 4, and 5, whereas the importance of the Pay Bursts Only Once phenomenon increases, if a larger CBS is used. In addition

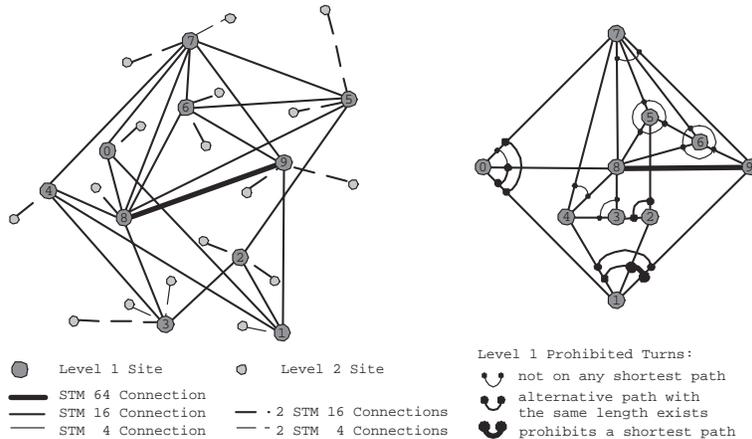


Fig. 3. DFN Topology and Example Level 1 Prohibited Turns

the Extended Pay Bursts Only Once principle is of major relevance, if the load on the network is high, and if large bursts are permitted. In case of a high load, flows are likely to share common sub-paths, and the interference of such flows with each other is much more accurately described by the Extended Pay Bursts Only Once principle. Thus, the form presented in (21), allows to derive tighter delay bounds, and thus to increase the acceptance ratio. In particular, as Figure 6 shows, the delay bound for the 99-percentile of the flows is reduced from above 112 ms to 81 ms in case of the Pay Bursts Only Once principle, and than to 59 ms in case of the extended principle for aggregate scheduling. For the 95-percentile, 85 ms, 63 ms, and 49 ms can be given. Further on, the load, up to which an acceptance ratio of for example 0.95 can be achieved, can be multiplied, when accounting for the Extended Pay Bursts Only Once phenomenon.

6 Conclusions

In this paper we have shown that a counterpart to the Pay Bursts Only Once phenomenon exists for interfering flows in aggregate scheduling networks. We then have derived a closed form solution for the end-to-end per-flow service curve in arbitrary feed forward aggregate scheduling networks, where links are of the rate-latency type, and flows are sigma-rho leaky bucket constraint. Our solution accounts for the known Pay Bursts Only Once principle, and extends it to aggregate scheduling in that bursts of interfering flows are paid only once, too. Thus, our form allows to give significantly closer bounds on the delay, while the intuitive form reduces computational complexity, if for example applied as a decision criterion for a Differentiated Services admission control in a Bandwidth Broker. A significant performance gain has been shown by simulation results.

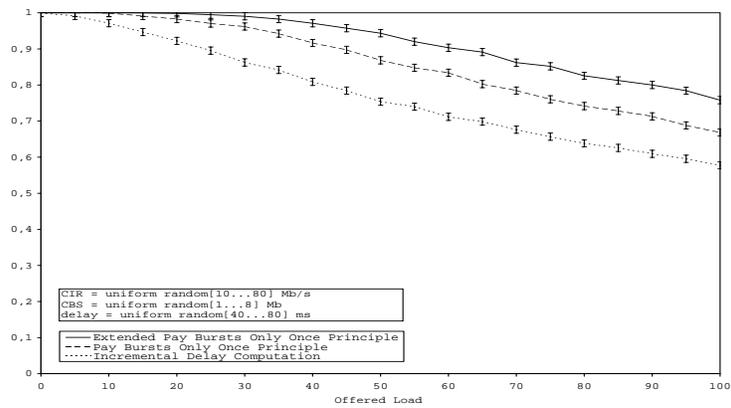


Fig. 4. Acceptance Ratio versus Load

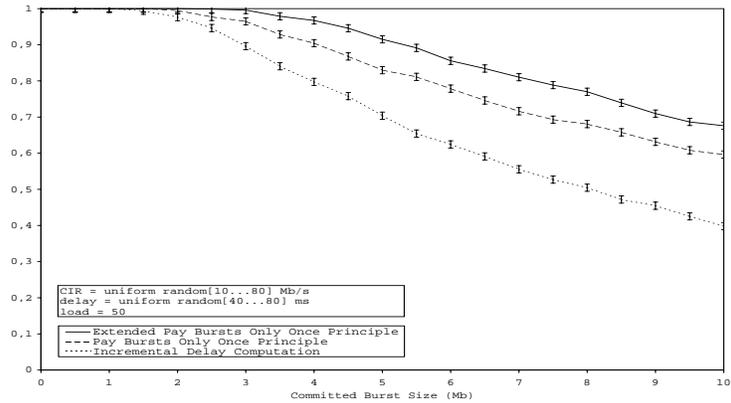


Fig. 5. Acceptance Ratio versus Committed Burst Size

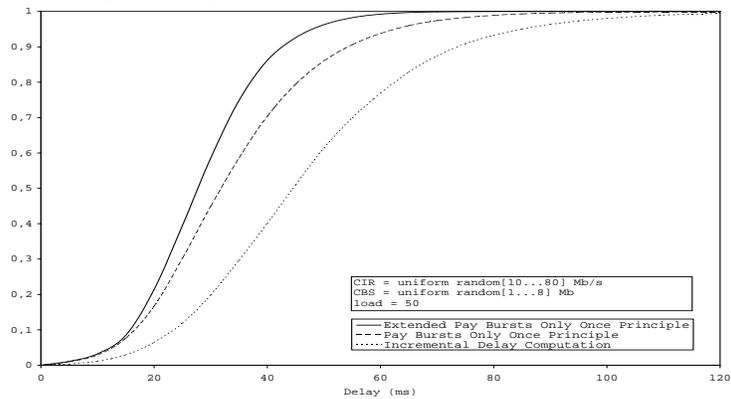


Fig. 6. Fraction of Flows with a smaller Delay Bound

Acknowledgements

This work was supported by the German Research Community (DFG) under grant graduate school (GRK) 643 "Software for Communication Systems".

References

1. Adler, H.-M., *10 Gigabit/s Plattform für das G-WiN betriebsbereit*, DFN Mitteilungen, Heft 60, November 2002.
2. Blake, S., et al., *An Architecture for Differentiated Services* RFC 2475, 1998.
3. Chang, C.-S., *Performance Guarantees in Communication Networks*, Springer, TNCS, 2000.
4. Charny, A., and Le Boudec, J.-Y., *Delay Bounds in a Network with Aggregate Scheduling*, Springer, LNCS 1922, Proceedings of QofIS, 2000.
5. Charny, A., et al., *Supplemental Information for the New Definition of EF PHB (Expedited Forwarding Per-Hop-Behavior)* RFC 3247, 2002.
6. Cruz, R. L., *A Calculus for Network Delay, Part I: Network Elements in Isolation*, IEEE Transactions on Information Theory, vol. 37, no. 1, pp 114-131, 1991.
7. Cruz, R. L., *A Calculus for Network Delay, Part II: Network Analysis*, IEEE Transactions on Information Theory, vol. 37, no. 1, pp 132-141, 1991.
8. Cruz, R. L., *SCED+: Efficient Management of Quality of Service Guarantees*, IEEE Infocom, 1998.
9. Davie, B., et al., *An Expedited Forwarding PHB* RFC 3246, 2002.
10. Foster, I., Fidler, M., Roy, A., Sander, V., Winkler, L., *End-to-End Quality of Service for High-End Applications* Elsevier Computer Communications Journal, in press, 2003.
11. Law, A. M., and Kelton, W. D., *Simulation, Modeling, and Analysis* McGraw-Hill, 3rd edition, 2000.
12. Le Boudec, J.-Y., and Thiran, P., *Network Calculus Made Easy*, Technical Report EPFL-DI 96/218. Ecole Polytechnique Federale, Lausanne (EPFL), 1996.
13. Le Boudec, J.-Y., and Thiran, P., *Network Calculus A Theory of Deterministic Queuing Systems for the Internet*, Springer, LNCS 2050, Version July 6, 2002.
14. Nichols, K., and Carpenter, B., *Definition of Differentiated Services Per Domain Behaviors and Rules for their Specification*, RFC 3086, 2001.
15. Nichols, K., Jacobson, V., and Zhang, L., *A Two-bit Differentiated Services Architecture for the Internet*, RFC 2638, 1999.
16. Pareck, A. K., and Gallager, R. G., *A Generalized Processor Sharing Approach to Flow Control in Integrated Services Networks: The Single-Node Case*, IEEE/ACM Transactions on Networking, vol. 1, no. 3, pp. 344-357, 1993.
17. Pareck, A. K., and Gallager, R. G., *A Generalized Processor Sharing Approach to Flow Control in Integrated Services Networks: The Multiple-Node Case*, IEEE/ACM Transactions on Networking, vol. 2, no. 2, pp. 137-150, 1994.
18. Paxson, V., and Floyd, S., *Wide Area Traffic: The Failure of Poisson Modeling* IEEE/ACM Transactions on Networking, vol. 3, no. 3, pp. 226-244, 1995.
19. Sander, V., *Design and Evaluation of a Bandwidth Broker that Provides Network Quality of Service for Grid Applications*, Ph.D. Thesis, Aachen University, 2002.
20. Sander, V., and Fidler, M., *Evaluation of a Differentiated Services based Implementation of a Premium and an Olympic Service*, Springer, LNCS 2511, Proceedings of QofIS, 2002.
21. Starobinski, D., Karpovsky, M., and Zakrevski, L., *Application of Network Calculus to General Topologies using Turn-Prohibition*, Proceedings of IEEE Infocom, 2002.