

# On the Benefits of Traffic “Reprofiling” The Single Hop Case

Jiayi Song, Jiaming Qiu, *Member, IEEE*, Roch Guérin, *Fellow, ACM, IEEE*, and Henry Sarioan, *Senior Member, IEEE*

**Abstract**—The need to guarantee hard delay bounds to traffic flows with deterministic traffic profiles, *e.g.*, token buckets, arises in several network settings. It is of interest to offer such guarantees while minimizing network bandwidth. The paper explores a basic building block, namely, a single hop configuration, towards realizing such a goal. The main results are in the form of optimal solutions for meeting local deadlines under schedulers of varying complexity and therefore cost. The results demonstrate how judiciously modifying flows’ traffic profiles, *i.e.*, *reprofiling* them, can help simple schedulers reduce the bandwidth they require, often performing nearly as well as more complex ones.

**Index Terms**—latency, bandwidth, optimization, token bucket, scheduling.

## I. INTRODUCTION

The provision of deterministic delay guarantees to traffic flows is emerging as an important requirement in increasingly diverse settings. They include automotive, avionics, and manufacturing applications, smart grids, and datacenters [1]–[7]. This is reflected in standards such as Time Sensitive Networking (TSN) and Deterministic Networking (DetNet) [8]–[11] and in the Service Level Objectives/Agreements (SLOs/SLAs) [12] of many service provider networks that are increasingly including latency targets, motivated in part by the rapid growth of edge computing offerings [13].

In such settings, the traffic eligible for latency guarantees is commonly controlled using a traffic regulator [14] in the form of a token bucket  $(r, b)$  that limits both the flow’s long-term rate,  $r$ , and burstiness,  $b$ . A flow’s token bucket parameters are typically determined using traces, and selected to ensure zero access delay [15]. The network’s goal is then to ensure that the latency guarantees of all such rate-controlled flows are met, preferably with as little bandwidth as possible.

This is the environment this paper assumes, with a focus on a basic building block, namely, delivering latency guarantees on a *single link* (hop) with *the least amount of bandwidth*. The answer obviously depends on the type of scheduler controlling access to the link, and the paper

considers schedulers of different levels of complexity. Of greater interest is whether, what the paper terms *reprofiling*, can be beneficial. Reprofile amounts to modifying a flow’s original (chosen by the user) token bucket parameters to make the flow “easier” to accommodate. This concept was explored in WorkloadCompactor [15] with one important difference, namely, the constraint that reprofiling should not introduce any delay. In contrast, our reprofiling solutions impose an added delay in exchange for smoother flows. This in turn calls for tighter network latency bounds to ensure that the original delay targets are still met. The outcome of this trade-off depends on the level of reprofiling applied as well as the type of scheduler in use. Investigating when and how it is positive in a single hop setting is the focus of this paper.

Specifically, the paper considers reprofiling of the form  $(r, b) \xrightarrow{\text{reprofiling}} (r, b')$ , where  $b' \leq b$ . In other words, we reduce the flow’s burstiness to make it easier to handle. We note that more complex reprofiling solutions are possible. Our motivations for focusing on burst reduction are two-fold. First, we want to minimize any added complexity, and this reprofiling can be realized simply by modifying the burst parameter of the *existing* token bucket. Second, As shown in [16, Appendix F], in simple configurations involving only two flows and a static priority scheduler, adjusting the burst size is sufficient to minimize the required bandwidth. These motivations notwithstanding, more complex reprofilers, *e.g.*, adding a second token bucket that controls the peak rate, can be of benefit in more general settings. We explore this extension in [17] in the multiple hops setting.

We note that the notion of reprofiling is closely tied to the definition of *greedy shapers* of [18, Section 1.5], with one important difference. Specifically, depending on the scheduler, a reprofiler can be either non-work-conserving, *i.e.*, as a (greedy) shaper, or work-conserving. The latter is only applicable when relying on dynamic priority schedulers such as earliest deadline first (edf) that can combine the local link deadline and the reprofiling delay when determining the order in which to send packets.

The paper makes the following contributions when it comes to meeting latency targets in the single-hop case with traffic profiles in the form of token buckets:

- Characterize the optimal (minimum bandwidth) solution, and show that a dynamic priority (edf) scheduler can realize it. The solution readily establishes that reprofiling yields no benefit with such a scheduler.
- Identify optimal reprofiling solutions for static priority

J. Qiu and R. Guérin are with the Computer Science and Engineering department at Washington University in St. Louis, Saint Louis, MO 63130, USA, e-mail: {qiujiaming, guerin}@wustl.edu.

J. Song was with the Computer Science and Engineering department at Washington University in St. Louis and is now with ByteDance Inc., San Jose, CA 95110, USA, e-mail: jiaiyisong@wustl.edu.

H. Sarioan is with Google, Mountain View, CA 94043, USA, e-mail: hsarioan@google.com.

An early version of the paper was presented at the 33rd International Teletraffic Congress (ITC’33) in September 2021.

The work was supported by NSF grant CNS 2006530 and a gift from Google.

and fifo schedulers, and demonstrate how they allow those schedulers to closely approximate the performance of the more complex edf scheduler across a range of scenarios.

For ease of exposition, the results are derived and presented assuming a fluid model, which, therefore, implies a preemptive behavior. As the discussion of [18, Section 1.1.1] highlights, extending the results to a packet setting is readily achievable from standard network calculus results. For illustration purposes, Appendix F of [16] derives a solution for a static priority scheduler under a packet-based model, but the results do not contribute further insight.

The paper is structured as follows. Section II introduces our traffic model and optimization framework. The next three sections present optimal solutions for schedulers of different complexity. Section III considers a general, dynamic priority scheduler, while Sections IV and V assume simpler static priority and fifo schedulers. For the latter two, the benefits of reprofiling flows are also explored. Section VI quantifies performance for each scheduler, starting with two-flow configurations that help build intuition for the results, before considering more general multi-flow scenarios. Section VII reviews related works, while Section VIII summarizes the paper's findings and their relevance to the multi-hop extension of [17]. Proofs and ancillary results are relegated to appendices available in an online version of the paper [16].

## II. MODEL FORMULATION

Consider the configuration of Fig. 1 with  $n$  flows sharing a common link<sup>1</sup> of rate  $R$ . The traffic generated by flow  $i$  is rate-controlled using a two-parameter token bucket  $(r_i, b_i)$  [14], its *traffic profile*, where  $r_i$  is the token rate and  $b_i$  the bucket size. Flow  $i$  also has a local packet-level deadline  $d_i$ , where w.l.o.g. we assume  $d_1 > d_2 > \dots > d_n$  with  $d_1 < \infty$ . Our goal is to meet the deadlines of all  $n$  flows with the lowest possible link bandwidth  $R$ . In doing so, we further assume *greedy sources* [18, Proposition 1.2.5] that fully realize the arrival curve associated with their token bucket.

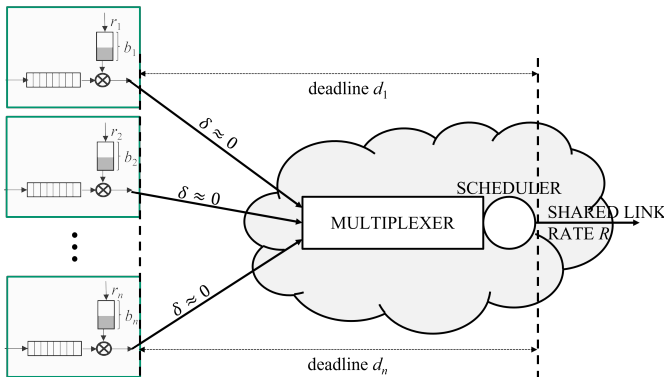


Fig. 1: A typical one-hop configuration with  $n$  flows.

In this setting, let  $\mathbf{r} = (r_1, r_2, \dots, r_n)$ ,  $\mathbf{b} = (b_1, b_2, \dots, b_n)$ , and  $\mathbf{d} = (d_1, d_2, \dots, d_n)$  be the vectors of rates, burst sizes, and deadlines of the flows sharing the link, respectively,

<sup>1</sup>For simplicity, we assume that enough buffering is available and that the link capacity is such that the system is stable and lossless.

and let  $D_i^*(\mathbf{r}, \mathbf{b}, R)$  denote flow  $i$ 's *worst-case* delay (queueing+transmission). Our bandwidth minimization problem can then be formulated as an optimization of the form:

$$\begin{aligned} \text{OPT-}\square \quad & \min R \\ \text{s.t} \quad & D_i^*(\mathbf{r}, \mathbf{b}, R) \leq d_i, \quad \forall i, 1 \leq i \leq n \end{aligned}$$

where  $R$  is the optimization variable and  $\square$  denotes the scheduler type, which for notational simplicity has been omitted in the expression of  $D_i^*(\mathbf{r}, \mathbf{b}, R)$ . As mentioned in Section I, one of our goals is to evaluate the trade-off between (bandwidth) efficiency and complexity across different schedulers.

Another goal is to investigate the potential benefits of *reprofiling* flows prior to forwarding their traffic to the scheduler. Reprofile amounts to applying a different, typically “smaller”, traffic profile to each flow before forwarding them to the scheduler. This can<sup>2</sup> introduce an up-front reprofiling delay, but may lower the bandwidth required to meet overall latency goals if it makes flows “easier” to handle.

More formally, given a scheduler “ $\square$ ” and  $n$  flows sharing a link, where flow  $i$ ,  $1 \leq i \leq n$ , has traffic profile  $(r_i, b_i)$  and deadline  $d_i$ , the goal of reprofiling is to identify smaller burst sizes  $b'_i \leq b_i$ ,  $1 \leq i \leq n$ , that minimize the link bandwidth  $R$  needed to meet the flows' deadlines, inclusive of any resulting reprofiling delay (the smaller burst  $b'_i$  introduces a reprofiling delay of  $\frac{b_i - b'_i}{r_i}$ ). We note that we restrict reprofiling options to only reducing the burst size, rather than also considering adding a “peak rate” shaper. This is in part to simplify the resulting optimization **OPT\_R** $\square$ , and also because, as shown in Appendix F of [16], this is sufficient in simple configurations with only two flows. This translates into a modified optimization problem **OPT\_R** $\square$  of the form

$$\begin{aligned} \text{OPT\_R}\square \quad & \min_{\mathbf{b}'} R \\ \text{s.t} \quad & D_i^*(\mathbf{r}, \mathbf{b}, \mathbf{b}', R) \leq d_i, \quad \forall i, 1 \leq i \leq n \end{aligned}$$

where  $R$  and  $\mathbf{b}'$  are the optimization variables. The latter denotes the vectors of updated (reprofiled) burst sizes of the  $n$  flows, and  $D_i^*(\mathbf{r}, \mathbf{b}, \mathbf{b}', R)$ ,  $1 \leq i \leq n$ , are the worst case delays, accounting for reprofiling delays, of the  $n$  flows under scheduler  $\square$  and a link bandwidth of  $R$ . The optimization explores the extent to which making flows smoother (smaller bursts) can facilitate meeting their delay targets with less bandwidth in spite of the access delay that reprofiling adds.

The next three sections explore solutions to **OPT** $\square$  (and **OPT\_R** $\square$ ) for different schedulers, namely, dynamic priority, static priority, and fifo (**OPT\_DP**, **OPT\_SP**, and **OPT\_F**).

## III. DYNAMIC PRIORITIES

We start with the most powerful but most complex scheduler, dynamic priorities, with priorities derived from service curves assigned to flows as a function of their profile (deadline and traffic envelope). We first solve **OPT\_DP** by characterizing the service curves that achieve the lowest bandwidth while meeting all deadlines in the absence of any reprofiling.

To derive the result, we first specify a service-curve assignment  $\Gamma_{sc}$  that satisfies all deadlines, identify the minimum

<sup>2</sup>When the reprofiler operates in a non-work-conserving manner.

link bandwidth  $R^*$  required to realize  $\Gamma_{sc}$ , and show that any scheduler requires at least  $R^*$ . We then show that an earliest deadline first scheduler realizes  $\Gamma_{sc}$  and, therefore, meets all the flow deadlines under  $R^*$ . We note that this then implies that reprofiling is of no benefit when an edf scheduler is available.

**Proposition 1.** *Consider a link shared by  $n$  token bucket controlled flows, where flow  $i, 1 \leq i \leq n$ , has a traffic profile  $(r_i, b_i)$  and a deadline  $d_i$ , with  $d_1 > d_2 > \dots > d_n$  and  $d_1 < \infty$ . Consider a service-curve assignment  $\Gamma_{sc}$  that allocates flow  $i$  a service curve of*

$$SC_i(t) = \begin{cases} 0 & \text{when } t < d_i, \\ b_i + r_i(t - d_i) & \text{otherwise.} \end{cases} \quad (1)$$

Then

- 1) For any flow  $i, 1 \leq i \leq n$ ,  $SC_i(t)$  ensures a worst-case end-to-end delay no larger than  $d_i$ .
- 2) Realizing  $\Gamma_{sc}$  requires a link bandwidth of at least

$$R^* = \max_{1 \leq h \leq n} \left\{ \sum_{i=1}^n r_i, \frac{\sum_{i=h}^n b_i + r_i(d_h - d_i)}{d_h} \right\}. \quad (2)$$

- 3) Any scheduling mechanism capable of meeting all the flows' deadlines requires a bandwidth of at least  $R^*$ .

The proof of Proposition 1 is in Appendix B-A of [16]. The optimality of  $\Gamma_{sc}$  is intuitive. Recall that a service curve is a lower bound on the service received by a flow. Eq. (1) assigns service to a flow at a rate exactly equal to its input rate, but delayed by its deadline, *i.e.*, provided at the latest possible time. Conversely, any mechanism  $\hat{\Gamma}$  that meets all flows' deadlines must by time  $t$  have provided flow  $i$  a cumulative service at least equal to the amount of data that flow  $i$  may have generated by time  $t - d_i$ , which is exactly  $SC_i(t)$ . Hence the mechanism must offer flow  $i$  a service curve  $\widehat{SC}_i(t) \geq SC_i(t), \forall t$ .

Next, we identify at least one mechanism capable of realizing the services curves of Eq. (1) under  $R^*$ , and consequently providing a solution to **OPT\_DP** for schedulers that support dynamic priorities.

**Proposition 2.** *Consider a link shared by  $n$  token bucket controlled flows, where flow  $i, 1 \leq i \leq n$ , has traffic profile  $(r_i, b_i)$  and deadline  $d_i$ , with  $d_1 > d_2 > \dots > d_n$  and  $d_1 < \infty$ . The earliest deadline first (edf) scheduler realizes  $\Gamma_{sc}$  under a link bandwidth of  $R^*$ .*

The proof of Proposition 2 is in Appendix B-B of [16]. We note that the optimality of edf is intuitive, as minimizing the required bandwidth is the dual problem to maximizing the schedulable region for which edf's optimality is known [19].

As previously mentioned and as the next proposition formally states, reprofiling does not reduce the minimum required bandwidth  $R^*$  of Eq. (2). Consequently, it affords no benefits with edf schedulers capable of meeting the deadlines under  $R^*$ . This is expected given the optimality of edf schedulers.

**Proposition 3.** *Consider a link shared by  $n$  token bucket controlled flows, where flow  $i, 1 \leq i \leq n$ , has traffic profile  $(r_i, b_i)$  and deadline  $d_i$ , with  $d_1 > d_2 > \dots > d_n$  and  $d_1 < \infty$ .*

*Reprofiling flows will not decrease the minimum bandwidth required to meet the flows' deadlines.*

The proof is in Appendix B-C of [16].

Note that  $\Gamma_{sc}$  specifies a non-linear (piece-wise-linear) service curve for each flow. Given the popularity and simplicity of linear service curves, *i.e.*, rate-based schedulers, it is tempting to investigate whether such schedulers, *e.g.*, GPS [20], could be used instead. Unfortunately, it is easy to find scenarios where linear service curves perform worse.

Consider a link shared by two flows with traffic profiles  $(r_1, b_1) = (1, 45)$  and  $(r_2, b_2) = (1, 5)$ , and deadlines  $d_1 = 10$  and  $d_2 = 1$ . A rate-based scheduler must allocate a bandwidth of  $\max\{\frac{b}{d}, r\}$  to a flow with traffic profile  $(r, b)$  to meet its deadline  $d$  of. Applying this to flow 2 that has the tighter deadline calls for a bandwidth of 5 to meet its deadline. After 1.25 units of time (the time to clear the initial burst of 5 and the additional data that accumulated during its transmission), flow 2's bandwidth usage drops down to  $r_2 = 1$ . The remaining 4 units then become available to flow 1. This means that the initial dedicated bandwidth needed by flow 1 to meet its deadline of 10 given its burst size of  $b_1 = 45$  is simply its token rate  $r_1 = 1^3$ , for a total network bandwidth of 6 units. In contrast, Eq. (2) tells us that  $\Gamma_{sc}$  only requires a bandwidth of  $R^* = 5.9$ .

The next two sections consider simpler static priority and fifo schedulers, and quantify the bandwidth they require to meet flows' deadlines. Both schedulers are considered either alone or with "reprofilers" that first modify the flows' traffic profiles before they are allowed to access the scheduler.

#### IV. STATIC PRIORITIES

Though edf schedulers are efficient and increasingly realizable [21]–[23], they are expensive and may not be practical in all environments. It is, therefore, of interest to explore simpler alternatives while quantifying the trade-off they entail between efficacy and complexity. For that purpose, we consider next a static priority scheduler where each flow is assigned a fixed priority as a function of its deadline.

As before, we consider  $n$  flows with traffic profiles  $(r_i, b_i)$  and deadlines  $d_i, 1 \leq i \leq n$ , sharing a common link. The question we first address is how to assign (static) priorities to each flow given their deadlines and **OPT\_SP**'s goal of minimizing link bandwidth? The next proposition offers a partial and somewhat intuitive answer to this question by establishing that the minimum link bandwidth can be achieved by giving flows with shorter deadlines a higher priority. Formally,

**Proposition 4.** *Consider a link shared by  $n$  token bucket controlled flows, where flow  $i, 1 \leq i \leq n$ , has traffic profile  $(r_i, b_i)$  and deadline  $d_i$ , with  $d_1 > d_2 > \dots > d_n$  and  $d_1 < \infty$ . Under a static-priority scheduler, there exists an assignment of flows to priorities that minimizes link bandwidth while meeting all flows deadlines such that flow  $i$  is assigned a priority strictly greater than that of flow  $j$  only if  $d_i < d_j$ .*

<sup>3</sup>Clearing the burst of flow 1 by its deadline  $d_1 = 10$  calls for a bandwidth  $x$  such that  $45 - \frac{5}{4}x - (x + 4)(10 - \frac{5}{4}) \leq 0$ , which yields  $x \geq 1$ .

The proof is in Appendix C-A of [16]. We note that while Proposition 4 states that link bandwidth can be minimized by assigning flows to priorities in the order of their deadline, it neither rules out other mappings nor does it imply that flows with different deadlines should always be mapped to distinct priorities. For example, large enough deadlines can all be met by a link bandwidth equal to the sum of the flows' average rates, *i.e.*,  $R^* = \sum_{i=1}^n r_i$ . In this case, priorities and their ordering are irrelevant. More generally, grouping flows with different deadlines in the same priority class can often result in a lower bandwidth than mapping them to distinct priority classes<sup>4</sup>. Nevertheless, motivated by Proposition 4, we propose a simple assignment rule that strictly maps lower deadline flows to higher priorities, and evaluate its performance.

#### A. Static Priorities without Reprofileing

From [18, Proposition 1.3.4] we know that when  $n$  flows with traffic profiles  $(r_i, b_i)$ ,  $1 \leq i \leq n$ , share a link of bandwidth  $R \geq \sum_{i=1}^n r_i$  with flow  $i$  assigned to priority  $i$  (priority  $n$  is the highest), then, under a static-priority scheduler, the worst case delay of flow  $h$  is upper-bounded by  $\frac{\sum_{i=h}^n b_i}{R - \sum_{i=h+1}^n r_i}$  (recall that under our notation, priority  $n$  is the highest). As a result, the minimum link bandwidth  $\tilde{R}^*$  to ensure that flow  $h$ 's deadline  $d_h$  is met for all  $h$ , *i.e.*, solving **OPT\_SP**, is given by:

$$\tilde{R}^* = \max_{1 \leq h \leq n} \left\{ \sum_{i=1}^n r_i, \frac{\sum_{i=h}^n b_i}{d_h} + \sum_{i=h+1}^n r_i \right\} \quad (3)$$

Towards evaluating the performance of a static priority scheduler versus that of an edf scheduler, we compare  $\tilde{R}^*$  with  $R^*$  through their relative difference, *i.e.*,  $\frac{\tilde{R}^* - R^*}{R^*}$ . For ease of comparison, we rewrite  $R^*$  as

$$R^* = \max_{1 \leq h \leq n} \left\{ \sum_{i=1}^n r_i, \frac{\sum_{i=h}^n b_i}{d_h} + \sum_{i=h+1}^n r_i \left( 1 - \frac{d_i}{d_h} \right) \right\} \quad (4)$$

Comparing Eqs. (3) and (4) shows that  $R^* = \tilde{R}^*$  iff  $\tilde{R}^* = \sum_{i=1}^n r_i$ , *i.e.*,  $\frac{\sum_{i=h}^n b_i}{d_h} \leq \sum_{i=h+1}^n r_i$ ,  $\forall 1 \leq h \leq n$ . In other words, static priority and edf schedulers perform equally well (yield the same minimum bandwidth), when flow bursts are small and deadlines relatively large so that they can be met with a link bandwidth equal to the sum of the token rates. However, when  $\tilde{R}^* \neq \sum_{i=1}^n r_i$ , a static priority scheduler can require a much larger bandwidth.

Consider a scenario where  $R^*$  is achieved at  $h^*$ , *i.e.*,  $R^* = \frac{\sum_{i=h^*}^n b_i}{d_{h^*}} + \sum_{i=h^*+1}^n r_i \left( 1 - \frac{d_i}{d_{h^*}} \right)$ . Though  $\tilde{R}^*$  may not be realized at the same  $h^*$  value, this still provides a lower bound for  $\tilde{R}^*$ , namely,  $\tilde{R}^* \geq \frac{\sum_{i=h^*}^n b_i}{d_{h^*}} + \sum_{i=h^*+1}^n r_i$ . Thus, the relative difference between  $\tilde{R}^*$  and  $R^*$  is no less than

$$\begin{aligned} & \frac{\frac{\sum_{i=h^*}^n b_i}{d_{h^*}} + \sum_{i=h^*+1}^n r_i}{\frac{\sum_{i=h^*}^n b_i}{d_{h^*}} + \sum_{i=h^*+1}^n r_i \left( 1 - \frac{d_i}{d_{h^*}} \right)} - 1 \\ &= \frac{\sum_{i=h^*+1}^n d_i r_i}{\sum_{i=h^*}^n b_i + \sum_{i=h^*+1}^n r_i (d_{h^*} - d_i)} \end{aligned} \quad (5)$$

As the right-hand-side of Eq. (5) increases with  $d_i$  for all  $i \geq h^*$ , it is maximized for  $d_i = d_{h^*} - \epsilon_i$ ,  $\forall i > h^*$ , for arbitrarily small  $\epsilon_{h^*+1} < \dots < \epsilon_n$ , so that its supremum is equal to  $\frac{\sum_{i=h^*+1}^n r_i d_{h^*}}{\sum_{i=h^*}^n b_i}$ . Note that this is intuitive, as when flows have arbitrarily close deadlines, they should receive equal service shares, which is in direct conflict with a strict priority ordering.

Under certain flow profiles, the above supremum can be large. In a two-flow scenario, basic algebraic manipulations give a supremum of  $\frac{r_2}{r_1 + r_2}$ , which is achieved at  $d_2 = d_1 = \frac{b_2 + b_1}{r_1 + r_2}$ . Since  $\frac{r_2}{r_1 + r_2} \rightarrow 1$  as  $\frac{r_1}{r_2} \rightarrow 0$ , the optimal static priority scheduler in the two-flow case could require twice as much bandwidth as the optimal edf scheduler.

#### B. Static Priorities with Reprofileing

Static priorities can require significantly more bandwidth than  $R^*$  mostly because they are a rather blunt instrument when it comes to fine-tuning the allocation of transmission opportunities as a function of packet deadlines. In particular, they often result in some flows experiencing a delay much lower than their target deadline.

This is intrinsic to the static structure of the scheduler and to our choice of an assignment that maps distinct deadlines to different priorities, but can be mitigated by anticipating and leveraging the “slack” in the delay of some flows. One such option is to use this slack towards reprofileing those flows, *i.e.*, make them “smoother”. Of interest then, is how to reprofile flows to maximize any resulting link bandwidth reduction?

Consider the trivial example of a single link shared by two flows with traffic profiles  $(r_1, b_1) = (1, 5)$  and  $(r_2, b_2) = (4, 5)$  and deadlines  $d_1 = 1.4, d_2 = 1.25$ . A strict static-priority scheduler requires a bandwidth  $\tilde{R}^* = 11.14$ . Assume next that we reprofile flow 2 to  $(r_2, b'_2) = (4, 0)$  before it enters the scheduler. The added reprofileing delay of  $(b_2 - b'_2)/r_2 = 1.25$  reduces the scheduling delay budget down to 0, but eliminates all burstiness. As a result, we only need a bandwidth of 7.57 (under a fluid model) to meet both flows' deadlines (a bandwidth of  $4 = r_2$  is still consumed by flow 2, but the remaining 3.57 is sufficient to allow flow 1 to meet its deadline). In other words, reprofileing flow 2 yields a bandwidth decrease of more than 30%. This simple example illustrates the benefits that judicious reprofileing can afford.

The next few propositions characterize the optimal reprofileing solution and the resulting bandwidth gains for a static priority scheduler and a set of flows and deadlines. We first derive expressions for flows' reprofileing and scheduling delays under static priorities, before obtaining the optimal reprofileing solution and the resulting minimum link bandwidth  $R_R^*$ .

Specifically, given  $n$  flows with initial traffic profiles  $(r_i, b_i)$ ,  $1 \leq i \leq n$ , deadlines  $d_1 > d_2 > \dots > d_n$ , a reprofileing solution  $(r_i, b'_i)$ ,  $1 \leq i \leq n$ , and a link of bandwidth  $R$ , Proposition 5 characterizes the worst case delay (reprofileing plus scheduling) of each flow, when a static priority scheduler assigns flow  $i$  priority  $i$  (shorter deadlines have higher priority). The result is used to formulate an optimization problem, **OPT\_RSP**, that seeks to minimize the link bandwidth required to meet individual flows' deadlines. The variables of the optimization are the reprofileing solution and the link bandwidth.

<sup>4</sup>We illustrate this in Appendix E of [16] for the case of two flows sharing a static priority scheduler.

Proposition 7 characterizes the minimum bandwidth  $\tilde{R}_R^*$  that **OPT\_RSP** can achieve, while Proposition 8 provides the optimal reprofiling solution.

Let  $\mathbf{b}' = (b'_1, b'_2, b'_3, \dots, b'_n)$  be the vector of reprofiled flow bursts, with  $B'_i = \sum_{j=i}^n b'_j$  and  $R_i = \sum_{j=i}^n r_j$ , the sum of the reprofiled bursts and rates of flows with priority greater than or equal to  $i$ ,  $1 \leq i \leq n$ , where  $B'_i = 0$  and  $R_i = 0$  for  $i > n$ . Flow  $i$ 's worst-case end-to-end delay is characterized next.

**Proposition 5.** Consider a link shared by  $n$  token bucket controlled flows, where flow  $i$ ,  $1 \leq i \leq n$ , has traffic profile  $(r_i, b_i)$ . Assume a static priority scheduler that assigns flow  $i$  a priority of  $i$ , where priority  $n$  is the highest priority, and reprofiles flow  $i$  to  $(r_i, b'_i)$ , where  $0 \leq b'_i \leq b_i$ . Given a link bandwidth of  $R \geq \sum_{j=1}^n r_j$ , the worst-case delay for flow  $i$  is

$$D_i^* = \max \left\{ \frac{b_i + B'_{i+1}}{R - R_{i+1}}, \frac{b_i - b'_i}{r_i} + \frac{B'_{i+1}}{R - R_{i+1}} \right\}. \quad (6)$$

The proof is in Appendix C-B of [16]. Note that Eq. (6) states that flow  $i$ 's worst-case delay is realized by the last bit of its burst. The two terms of Eq. (6) capture the cases when this bit arrives before or after the end of flow  $i$ 's last busy period at the link, respectively, as this determines the extent to which it is affected by the reprofiling delay.

Observe also that  $D_i^*$  is independent of  $b'_1$  for  $2 \leq i \leq n$ , and decreases with  $b'_1$  when  $i = 1$ . This is intuitive as flow 1 has the lowest priority so that reprofiling it can neither decrease the worst-case end-to-end delay of other flows, nor consequently reduce the minimum link bandwidth required to meet specific deadlines for each flow. Formally,

**Corollary 6.** Consider a link shared by  $n$  token bucket controlled flows, where flow  $i$ ,  $1 \leq i \leq n$ , has traffic profile  $(r_i, b_i)$  and deadline  $d_i$ , with  $d_1 > d_2 > \dots > d_n$  and  $d_1 < \infty$ . Assume a static priority scheduler that assigns flow  $i$  a priority of  $i$ , where priority  $n$  is the highest priority, and reprofiles flow  $i$  to  $(r_i, b'_i)$ , where  $0 \leq b'_i \leq b_i$ . Given a link bandwidth of  $R \geq \sum_{j=1}^n r_j$ , reprofiling flow 1 cannot reduce the minimum required bandwidth.

Combining Proposition 5 and Corollary 6 with **OPT\_SP** gives the following optimization **OPT\_RSP** for a link shared by  $n$  flows and relying on a static priority scheduler preceded by reprofiling. Note that since the minimum link bandwidth needs to satisfy  $R \geq \sum_{i=1}^n r_i$ , combining this condition with  $R_i$ 's definition gives  $\sum_{i=1}^n r_i = R_1 \leq R$ .

$$\begin{aligned} \mathbf{OPT\_RSP} \quad & \min_{\mathbf{b}'} R \quad \text{s.t.} \\ & \max \left\{ \frac{b_i + B'_{i+1}}{R - R_{i+1}}, \frac{b_i - b'_i}{r_i} + \frac{B'_{i+1}}{R - R_{i+1}} \right\} \leq d_i, \quad \forall 1 \leq i \leq n, \\ & R_1 \leq R, \quad b'_1 = b_1, \quad 0 \leq b'_i \leq b_i, \quad \forall 2 \leq i \leq n. \end{aligned}$$

The solution of **OPT\_RSP** is characterized in Propositions 7 and 8 whose proofs are in Appendix C-C of [16]. Proposition 7 gives the optimal bandwidth  $\tilde{R}_R^*$  based only on flow profiles, and while it is too complex to yield a closed-form expression, it offers a feasible numerical procedure to compute  $\tilde{R}_R^*$ .

**Proposition 7.** For  $1 \leq i \leq n$ , denote  $H_i = b_i - d_i r_i$ ,  $\Pi_i(R) = \frac{r_i + R - R_{i+1}}{R - R_{i+1}}$  and  $V_i(R) = d_i(R - R_{i+1}) - b_i$ . Define  $\mathbb{S}_1(R) = \{V_1(R)\}$ , and  $\mathbb{S}_i(R) = \mathbb{S}_{i-1}(R) \cup \{V_i(R)\} \cup \left\{ \frac{s - H_i}{\Pi_i(R)} \mid s \in \mathbb{S}_{i-1}(R) \right\}$  for  $2 \leq i \leq n$ . Then we have  $\tilde{R}_R^* = \max \{R_1, \inf \{R \mid \forall s \in \mathbb{S}_n(R), s \geq 0\}\}$ .

Computing  $\tilde{R}_R^*$  requires solving polynomial inequalities of degree  $(n-1)$ , so that a closed-form expression is not feasible except for small  $n$ . However, as  $\mathbb{S}_i(R)$  relies only on flow profiles and  $\mathbb{S}_j(R)$ ,  $\forall j < i$ , we can recursively construct  $\mathbb{S}_n(R)$  from  $\mathbb{S}_1(R)$ . Hence, since  $R_1 \leq \tilde{R}_R^* \leq \tilde{R}^*$ , we can use a binary search to compute  $\tilde{R}_R^*$  from the relation  $\tilde{R}_R^* = \max \{R_1, \inf \{R \mid \forall s \in \mathbb{S}_n(R), s \geq 0\}\}$  in Proposition 7.

Next, Proposition 8 gives a constructive procedure to obtain the optimal reprofiling burst sizes  $\mathbf{b}'^*$  given  $\tilde{R}_R^*$  and the original flow profiles.

**Proposition 8.** The optimal reprofiling solution  $\mathbf{b}'^*$  satisfies

$$b_i'^* = \begin{cases} \max\{0, b_n - r_n d_n\}, & i = n; \\ \max \left\{ 0, b_i - r_i d_i + \frac{r_i B'_{i+1}^*}{\tilde{R}_R^* - R_{i+1}} \right\}, & 2 \leq i \leq n-1. \end{cases} \quad (7)$$

where we recall that  $b_1^* = b_1$  and  $B_i^* = \sum_{j=i}^n b_j^*$ .

Note that the optimal reprofiling burst size  $b_i'^*$  of flow  $i$ ,  $1 < i < n$  relies only on the optimal link bandwidth  $\tilde{R}_R^*$  and the reprofiling burst sizes of higher priority flows. Hence, we can recursively characterize  $b_i'^*$  from  $b_n^*$  given  $\tilde{R}_R^*$ .

## V. BASIC FIFO WITH REPROFILING

In this section, we consider a simple first-in-first-out (fifo) scheduler that serves data in the order in which it arrives. For conciseness and given the benefits of reprofiling demonstrated in Section IV-B, we directly assume that flows are reprofiled prior to being scheduled. Considering again a link shared by  $n$  flows with traffic profiles  $(r_i, b_i)$ ,  $1 \leq i \leq n$ , and deadlines  $d_1 > d_2 > \dots > d_n$ , our goal is to find a reprofiling solution  $(r_i, b'_i)$ ,  $1 \leq i \leq n$ , to minimize the link bandwidth required to meet the flows' deadlines.

Towards answering this question, we first proceed to characterize the worst case delay across  $n$  flows sharing a link of bandwidth  $R$  equipped with a fifo scheduler when the flows have initial traffic profiles  $(r_i, b_i)$ ,  $1 \leq i \leq n$ , and are reprofiled to  $(r_i, b'_i)$ ,  $1 \leq i \leq n$ , prior to being scheduled. Using this result, we then identify the reprofiled burst sizes  $b'_i$ ,  $1 \leq i \leq n$ , that minimize the link bandwidth required to ensure that all deadlines  $d_1 > d_2 > \dots > d_n$ , and  $d_1 < \infty$  are met. As with other configurations, we only state the results with proofs relegated to Appendix D of [16].

**Proposition 9.** Consider a system with  $n$  token bucket controlled flows with traffic profiles  $(r_i, b_i)$ ,  $1 \leq i \leq n$ , sharing a fifo link with bandwidth  $R \geq R_1 = \sum_{j=1}^n r_j$ . Assume that the system reprofiles flow  $i$  to  $(r_i, b'_i)$ . The worst-case delay for flow  $i$  is then

$$\hat{D}_i^* = \max \left\{ \frac{b_i - b'_i}{r_i} + \frac{\sum_{j \neq i} b'_j}{R}, \frac{\sum_{j=1}^n b'_j}{R} + \frac{(b_i - b'_i)R_1}{r_i R} \right\}. \quad (8)$$

The proof of Proposition 9 is in Appendix D-A of [16].

With the result of Proposition 9 in hand, we can formulate a corresponding optimization problem, **OPT\_RF**, for computing the optimal reprofiling solution that minimizes the link bandwidth required to meet the deadlines  $d_1 > d_2 > \dots > d_n$ , and  $d_1 < \infty$  of the  $n$  flows. Specifically, combining Proposition 9 with **OPT\_F** gives the following optimization **OPT\_RF** for a link shared by  $n$  flows when relying on a fifo scheduler preceded by reprofiling. As before,  $\sum_{i=1}^n r_i = R_1 \leq R$ .

$$\begin{aligned} \mathbf{OPT\_RF} \quad & \min_{b'} R \quad \text{s.t.} \quad \forall 1 \leq i \leq n \\ & \max \left\{ \frac{b_i - b'_i}{r_i} + \frac{\sum_{j \neq i} b'_j}{R}, \frac{\sum_{j=1}^n b'_j}{R} + \frac{(b_i - b'_i)R_1}{r_i R} \right\} \leq d_i, \\ & R_1 \leq R, \quad 0 \leq b'_i \leq b_i, \forall 1 \leq i \leq n. \end{aligned} \quad (9)$$

The solution of **OPT\_RF** is characterized in Propositions 10 and 11 with proofs in Appendix D-B of [16]. As with a static priority scheduler, Proposition 10 gives a numerical procedure to compute the optimal bandwidth  $\hat{R}_R^*$  given the flows' profiles, while Proposition 11 gives the optimal reprofiling solution  $\hat{b}^*$  given  $\hat{R}_R^*$  and the original flows' profiles.

**Proposition 10.** For  $1 \leq i \leq n$ , define  $H_i = b_i - d_i r_i$ ,  $\hat{B}_i = \sum_{j=1}^i b_j$ , and  $\mathbb{Z}_i = \{1 \leq j \leq i \mid j \in \mathbb{Z}\}$ . Denote

$$X_F(R) = \max_{\substack{P_1, P_2 \subseteq \mathbb{Z}_n, \\ P_2 \neq \emptyset, P_1 \cap P_2 = \emptyset}} \frac{\sum_{i \in P_1} \frac{RH_i}{R+r_i} + \sum_{i \in P_2} \left( b_i - \frac{r_i d_i R}{R_1} \right)}{1 - \sum_{i \in P_1} \frac{r_i}{R+r_i} - \sum_{i \in P_2} \frac{r_i}{R_1}}$$

and

$$Y_F(R) = \min_{1 \leq i \leq n-1} \left\{ \hat{B}_n, R d_n, \min_{\substack{P_1, P_2 \subseteq \mathbb{Z}_i, \\ P_1 \cap P_2 = \emptyset, \\ P_1 \cup P_2 \neq \emptyset}} \left\{ \frac{\hat{B}_i - \sum_{j \in P_1} \frac{RH_j}{R+r_j} - \sum_{j \in P_2} \left( b_j - \frac{r_j d_j R}{R_1} \right)}{\sum_{j \in P_1} \frac{r_j}{R+r_j} + \sum_{j \in P_2} \frac{r_j}{R_1}} \right\} \right\}.$$

Then the optimal solution for **OPT\_RF** is

$$\hat{R}_R^* = \max \left\{ R_1, \frac{\hat{B}_n R_1}{\sum_{i=1}^n r_i d_i}, \min \{ R \mid X_F(R) \leq Y_F(R) \} \right\}.$$

As  $\max \left\{ R_1, \frac{\hat{B}_n R_1}{\sum_{i=1}^n r_i d_i} \right\} \leq \hat{R}_R^* \leq \hat{R}^* = \max \left\{ R_1, \frac{\hat{B}_n}{d_n} \right\}$ , where  $\hat{R}^*$  is the minimum required bandwidth achieved by a base (no reprofiling) fifo system, we can use Proposition 10 and a binary search to compute  $\hat{R}_R^*$ . Once  $\hat{R}_R^*$  is known, Proposition 11 gives the optimal reprofiling solution.

**Proposition 11.** For  $1 \leq i \leq n$ , define  $T_i(\hat{B}'_n, R) = \max \left\{ 0, \frac{R}{R+r_i} \left( H_i + \frac{r_i}{R} \hat{B}'_n \right), b_i + \frac{r_i (\hat{B}'_n - R d_i)}{R_1} \right\}$ . The optimal reprofiling solution  $\hat{b}^*$  of **OPT\_RF**'s is given by  $\hat{b}_1^* = \hat{B}_1^*$ , and  $\hat{b}_i^* = \hat{B}_i^* - \hat{B}_{i-1}^*$ ,  $2 \leq i \leq n$ , where  $\hat{B}_i^*$  satisfy

$$\begin{cases} \hat{B}_n^* = X_F(\hat{R}_R^*), \\ \hat{B}_i^* = \max \left\{ \sum_{j=1}^i T_j(\hat{B}_n^*, \hat{R}_R^*), \hat{B}_{i+1}^* - b_{i+1} \right\}, \\ \text{when } 1 \leq i \leq n-1. \end{cases} \quad (10)$$

Note that  $\hat{B}_n^*$  relies only on  $\hat{R}_R^*$  and flows' profiles. Whereas when  $1 \leq i \leq n-1$ ,  $\hat{B}_i^*$  relies only on  $\hat{R}_R^*$ ,  $\hat{B}_n^*$ ,  $\hat{B}_{i+1}^*$  and flows' profiles. Hence, we can recursively characterize  $\hat{B}_i^*$  from  $\hat{B}_n^*$  given  $\hat{R}_R^*$ .

## VI. EVALUATION

In this section, we explore the relative benefits of the solutions developed in the previous three sections. Of interest is assessing the “cost of simplicity,” namely, the amount of additional bandwidth required by simpler schedulers such as static priority or fifo compared to an edf scheduler. Also of interest is the magnitude of the improvements that reprofiling affords with static priority and fifo schedulers. To that end, the evaluation proceeds with a number of *pairwise comparisons* to quantify the relative (bandwidth) cost of each alternative.

The evaluation first focuses (Section VI-A) on scenarios with just two flows. Closed-form expressions are then available for the minimum bandwidth of each configuration, which make formal comparisons possible. Section VI-B extends this to more “general” scenarios involving multiple flows with different combinations of deadlines and traffic profiles.

In the initial two-flow comparisons of Section VI-A, we first select a pair of representative traffic profiles (token buckets), and then vary the flows' respective deadlines over a wide range of values. For each such combination, we explicitly compute the relative differences in bandwidth required by the different schedulers (with and without reprofiling, as applicable) using expressions derived from the propositions obtained in the previous sections. The results are presented in the form of “heat-maps” across the range of deadline combinations.

For the more general scenarios involving multiple flows (Section VI-B), we first generate a set of flow profiles, *i.e.*, token buckets and deadlines, by randomly selecting them from within specified ranges. For each such combination, the amount of bandwidth required to meet the flows' deadlines are then computed using again results from the propositions derived in the previous sections. Finally, for each pair of schedulers, we report statistics (means, standard deviations and the 95% confidence intervals of the means) of the relative bandwidth differences across those random selections.

### A. Basic Two-Flow Configurations

Recalling our earlier notation for the minimum bandwidth in each configuration, *i.e.*,  $R^*$  (edf);  $\tilde{R}^*$  (static priority);  $\hat{R}_R^*$  (static priority w/ reprofiling);  $\hat{R}^*$  (fifo); and  $\hat{R}_R^*$  (fifo w/ reprofiling), and specializing Eq. (2) to a configuration with two flows,  $(r_1, b_1)$  and  $(r_2, b_2)$ , the absolute minimum bandwidth to meet the flows' deadlines  $d_1$  and  $d_2$  is given by

$$R^* = \max \left\{ r_1 + r_2, \frac{b_2}{d_2}, \frac{b_1 + b_2 - r_2 d_2}{d_1} + r_2 \right\}, \quad (11)$$

which is then also the bandwidth required by the edf scheduler.

Similarly, if we consider a static priority scheduler, from Eq. (3), its bandwidth requirement  $\tilde{R}^*$  (in the absence of any reprofiling) for the same two-flow configuration is of the form

$$\tilde{R}^* = \max \left\{ r_1 + r_2, \frac{b_2}{d_2}, \frac{b_1 + b_2}{d_1} + r_2 \right\}; \quad (12)$$

If (optimal) reprofiling is introduced, specializing Proposition 7 to two flows, the minimum bandwidth  $\hat{R}_R^*$  reduces to

$$\begin{aligned} & \max \left\{ r_1 + r_2, \frac{b_2}{d_2}, \frac{b_1 + b_2 - r_2 d_2}{d_1} + r_2 \right\}, \text{ when } \frac{b_2}{r_2} \geq \frac{b_1}{r_1} \\ & \max \left\{ r_1 + r_2, \frac{b_2}{d_2}, \frac{b_1 + \max \{b_2 - r_2 d_2, 0\}}{d_1} + r_2 \right\}, \\ & \text{otherwise;} \end{aligned} \quad (13)$$

Finally, specializing the results of Propositions 10 and 11 to two flows, we find that the minimum required bandwidth  $\hat{R}^*$  under fifo without reprofiling is

$$\hat{R}^* = \max \left\{ r_1 + r_2, \frac{b_1 + b_2}{d_2} \right\}; \quad (14)$$

and that when (optimal) reprofiling is used,  $\hat{R}_R^*$  is given by Eq. (15). With these expressions in hand, we can now assess the relative benefits of each option in this two-flow scenario.

Specifically, we consider next combinations consisting of two flows with representative token bucket parameters  $(r_1, b_1) = (4, 10)$  and  $(r_2, b_2) = (10, 18)$ , and systematically vary their respective deadlines  $d_1 \geq d_2$  over a range of values. The bandwidth required to meet the deadlines is then compared for different pairs of schedulers using the expressions reported in Eq. (11), Eq. (13), and Eq. (15).

1) *The Impact of Scheduler Complexity:* We first evaluate the impact of relying on schedulers of decreasing complexity, when those schedulers are coupled with an optimal reprofiling solution. In other words, we compare the bandwidth requirements of an edf scheduler to those of static priority and fifo schedulers combined with an optimal reprofiler. The comparison is in the form of relative differences (improvements realizable from more complex schedulers), i.e.,  $\frac{\hat{R}_R^* - R^*}{\hat{R}_R^*}$ ,  $\frac{\hat{R}_R^* - R^*}{\hat{R}_R^*}$ , and  $\frac{\hat{R}_R^* - \hat{R}_R^*}{\hat{R}_R^*}$ .

*edf vs. static priority w/ optimal reprofiling.*

We start with comparing an edf scheduler with a static priority scheduler plus optimal reprofiling. Eqs. (11) and (13) then state that  $R^* < \hat{R}_R^*$  iff  $\frac{b_2}{r_2} < d_2 \leq d_1 < \frac{b_1}{r_1}$ .

The results are reported in Fig. 2a, and, as mentioned, are in the form of a heat-map of the relative bandwidth differences as the flows' respective deadlines vary. As shown in the figure, a static priority scheduler, when combined with reprofiling, performs as well as an edf scheduler, except for a relatively small (triangular) region where  $d_1$  and  $d_2$  are close to each other and both of intermediate values<sup>5</sup>. Towards better characterizing this range, i.e.,  $d_2 > \frac{b_2}{r_2}$  and  $d_1 < \frac{b_1}{r_1}$ , we see that the supremum of  $\frac{\hat{R}_R^* - R^*}{\hat{R}_R^*}$  is achieved at  $d_1 = d_2 = \frac{b_1 + b_2}{r_1 + r_2}$ ,

with  $\hat{R}_R^* = \frac{b_1}{d_1} + r_2$ , and  $R^* = r_1 + r_2$ . The relative difference in bandwidth between the two schemes is then of the form

$$\frac{\hat{R}_R^* - R^*}{\hat{R}_R^*} = 1 - \frac{1}{\frac{b_1}{b_1 + b_2} + \frac{r_2}{r_1 + r_2}},$$

which can be shown to be upper-bounded by 0.5. In other words, in the two-flow case, the (optimal) edf scheduler can result in a bandwidth saving of at most 50% when compared to a static priority scheduler with (optimal) reprofiling. This happens when the deadlines of the two flows are very close to each other, a scenario unlikely in practice.

*edf vs. fifo w/ optimal reprofiling*

Next, we compare an edf scheduler and a fifo scheduler plus optimal reprofiling. Eqs. (11) and (15) state that  $\hat{R}_R^* > R^*$  iff  $d_1 - \frac{b_1}{r_1} < d_2 < \frac{b_1 + b_2 - d_1 r_1}{r_2}$ . We illustrate the corresponding relative differences in Fig. 2b using the same two-flow combination as before. From the figure, we see that fifo + reprofiling performs poorly relative to an edf scheduler when neither  $d_1$  nor  $d_2$  are large. As with static priorities, such configurations may not be common in practice.

We note that the supremum of  $\frac{\hat{R}_R^* - R^*}{\hat{R}_R^*}$  is achieved when  $0 < d_2 < \frac{b_1 + b_2 + r_2 d_1 - \sqrt{(b_1 + b_2 + r_2 d_1)^2 - 4 r_2 b_2 d_1}}{2 r_2}$ , with Eq. (11) defaulting to  $R^* = \frac{b_2}{d_2}$  and Eq. (15) to  $\hat{R}_R^* = \frac{b_1 + b_2 - d_1 r_1 + \sqrt{(b_1 + b_2 - d_1 r_1)^2 + 4 r_1 d_2 b_2}}{2 d_2}$ . Hence, the relative difference becomes

$$\begin{aligned} \frac{\hat{R}_R^* - R^*}{\hat{R}_R^*} &= 1 \\ &- \frac{2 b_2}{b_1 + b_2 - d_1 r_1 + \sqrt{(b_1 + b_2 - d_1 r_1)^2 + 4 r_1 d_2 b_2}}, \end{aligned}$$

which increases with  $d_2$ . Thus, its supremum is achieved as  $d_2 \rightarrow d_1$ . Similarly, one easily shows that  $1 - \frac{2 b_2}{b_1 + b_2 - d_1 r_1 + \sqrt{(b_1 + b_2 - d_1 r_1)^2 + 4 r_1 d_2 b_2}}$  decreases with  $d_1$ . Hence, the supremum of the relative difference is achieved as  $d_1 \rightarrow 0$ , and is of the form  $\frac{b_1}{b_1 + b_2}$ , which goes to 1 as  $\frac{b_1}{b_2} \rightarrow \infty$ . In other words, an edf scheduler can yield a 100% improvement over a fifo scheduler with optimal reprofiling.

*Fifo vs. static priority both w/ optimal reprofiling*

Finally, we compare fifo and static priority schedulers when both rely on optimal reprofiling. Eqs. (13) and (15) give that  $\hat{R}_R^* > \tilde{R}_R^*$  iff  $\max \left\{ \frac{b_2}{r_2}, \frac{(b_1 + b_2)(r_1 + r_2)}{r_2(b_1/d_1 + r_2)} \right\} < d_1 < \frac{b_1}{r_1}$ . Fig. 2c illustrates the difference, again relying on a heat-map for the same two-flow combination as the two previous scenarios.

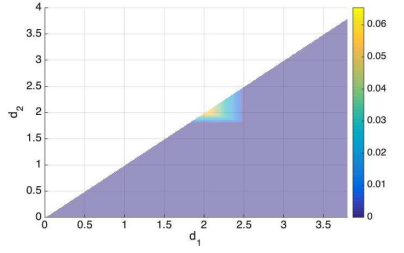
The figure shows that the benefits of priority are maximum when  $d_2$  is small and  $d_1$  is not too large. This is intuitive in that a small  $d_2$  calls for affording maximum protection to flow 2, which a priority structure offers more readily than a fifo. Conversely, when  $d_1$  is large, flow 1 can be reprofiled to eliminate all burstiness, which limits its impact on flow 2 even when both flows compete in a fifo scheduler.

The figure also reveals that a small region exists (when  $d_1$  and  $d_2$  are close to each other and both are of intermediate value) where fifo outperforms static priority. As alluded to

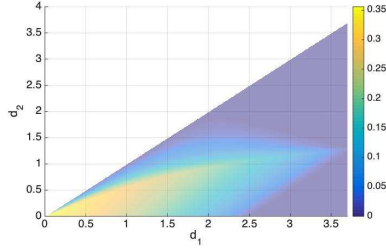
<sup>5</sup>(i) When  $d_2$  and  $d_1$  are close and small, the bandwidth required to meet the deadlines is very large under either edf or static priority schedulers. This ensures that both produce similar transmissions' orders. Consider, for example, a low-priority (larger deadline) burst that arrives  $(d_1 - d_2)$  before a high-priority (smaller deadline) one. It has higher priority under edf, and the speed of the link ensures it is transmitted before the arrival of the high-priority burst, which ensures no difference between edf and a static priority scheduler.  
(ii) When  $d_2$  and  $d_1$  are close but large, both schedulers meet their deadlines with the same bandwidth, i.e., the sum of the flows' average rates.



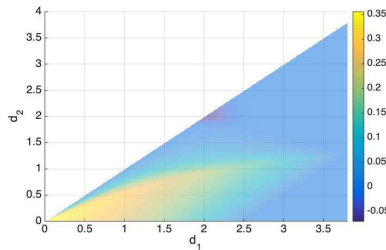
$$\hat{R}_R^* = \max \left\{ r_1 + r_2, \frac{b_2}{d_2}, \frac{(b_1 + b_2)(r_1 + r_2)}{d_1 r_1 + d_2 r_2}, \frac{b_1 + b_2 - d_1 r_1 + \sqrt{(b_1 + b_2 - d_1 r_1)^2 + 4r_1 d_2 b_2}}{2d_2} \right\}. \quad (15)$$



(a) Dyn. prio. vs. stat. prio. + reprofiling



(b) Dyn. prio. vs. fifo + reprofiling



(c) fifo + reprofiling vs. stat. prio. + reprofiling

Fig. 2: Relative bandwidth increases for  $(r_1, b_1) = (4, 10)$  and  $(r_2, b_2) = (10, 18)$ , as a function of  $d_1$  and  $d_2 < d_1$ . The figure is in the form of a heat-map. Darker colors (purple) correspond to smaller increases than lighter ones (yellow).

in the discussion following Proposition 4 and as expanded in Appendix E of [16], this is because a *strict* priority ordering of flows as a function of their deadlines needs not always be optimal. For instance, it is easy to see that two otherwise identical flows that only differ infinitesimally in their deadlines should be treated “identically.” This is more readily accomplished by having them share a common fifo queue than assigned to two distinct priorities.

To better understand differences in performance between the two schemes, we characterize the supremum and the infimum of  $\frac{\hat{R}_R^* - \tilde{R}_R^*}{\hat{R}_R^*}$ . Basic algebraic manipulations show that the supremum is achieved as  $d_1 = d_2 \rightarrow 0$ , where Eq. (15) defaults to  $\hat{R}_R^* = \frac{b_1 + b_2 - d_1 r_1 + \sqrt{(b_1 + b_2 - d_1 r_1)^2 + 4r_1 d_2 b_2}}{2d_2}$  and Eq. (13) to  $\tilde{R}_R^* = \frac{b_2}{d_2}$ , so that their relative difference is

ultimately of the form

$$\frac{\hat{R}_R^* - \tilde{R}_R^*}{\hat{R}_R^*} = \frac{b_1}{b_1 + b_2},$$

which goes to 1 as  $\frac{b_1}{b_2} \rightarrow \infty$ , *i.e.*, a maximum penalty of 100% for fifo with reprofiling over static priorities with reprofiling.

Conversely, the infimum is achieved at  $d_1 = d_2 = \frac{b_1 + b_2}{r_1 + r_2}$ , with Eqs. (13) and (15) defaulting to  $\tilde{R}_R^* = \frac{b_1}{d_1} + r_2$  and  $\hat{R}_R^* = r_1 + r_2$ , and a relative difference of the form

$$\frac{\hat{R}_R^* - \tilde{R}_R^*}{\hat{R}_R^*} = \frac{r_1}{r_1 + r_2} - \frac{b_1}{b_1 + b_2},$$

which increases with  $\frac{r_1}{r_2}$  and decreases with  $\frac{b_1}{b_2}$ . When  $\frac{r_1}{r_2} \rightarrow 0$  and  $\frac{b_1}{b_2} \rightarrow \infty$ , it achieves an infimum of  $-1$ , *i.e.*, a maximum penalty of 100% but now for static priorities with reprofiling over fifo with reprofiling. In other words, when used with reprofiling, both fifo and static priority can require twice as much bandwidth as the other.

Ensuring that static priority always outperforms fifo calls for determining when flows should be grouped in the same priority class rather than assigned to separate classes. Such grouping can be identified in simple scenarios with two or three flows, *e.g.*, see Appendix E of [16], but a general solution appears challenging. However, as we shall see in Section VI-B, the simple strict priority assignment on which we rely performs well in practice across a broad range of flow configurations.

2) *The Benefits of Reprofileing*: In this section, we evaluate the benefits afforded by (optimally) reprofileing flows with static priority and fifo schedulers. This is done by computing for both schedulers the minimum bandwidth required to meet flows’ deadlines without and with reprofileing, and evaluating the resulting relative differences, *i.e.*,  $\frac{\hat{R}_R^* - \tilde{R}_R^*}{\hat{R}_R^*}$  and  $\frac{\hat{R}_R^* - \tilde{R}_R^*}{\tilde{R}_R^*}$ .

For a static priority scheduler, Eqs. (12) and (13) indicate that  $\tilde{R}_R^* < \hat{R}_R^*$  iff  $\tilde{R}_R^* = \frac{b_1 + b_2}{d_1} + r_2 > \max \left\{ r_1 + r_2, \frac{b_2}{d_2} \right\}$ , *i.e.*, for a static priority scheduler, reprofileing<sup>6</sup> decreases the required bandwidth only when  $d_1$ , the larger deadline, is not too large and  $d_2$ , the smaller deadline, is not too small. This is intuitive. When  $d_1$  is large, the low-priority flow 1 can meet its deadline even without any mitigation of the impact of flow 2. Conversely, a small  $d_2$  offers little to no opportunity for reprofileing flow 2 as the added delay it introduces would need to be compensated by an even higher link bandwidth. This is illustrated in Fig. 3a for the same two-flow combination as in Fig. 2, *i.e.*,  $(r_1, b_1) = (4, 10)$  and  $(r_2, b_2) = (10, 18)$ . The intermediate region where “ $d_1$  is not too large and  $d_2$  is not too small” corresponds to the yellow triangular region where the benefits of reprofileing can reach 40%.

Similarly, Eqs. (14) and (15) indicate that  $\hat{R}_R^* < \tilde{R}_R^*$  iff  $d_2 < \frac{b_1 + b_2}{r_1 + r_2}$ , *i.e.*, for a fifo scheduler, reprofileing decreases

<sup>6</sup>Recall from Corollary 6 that the flow with the largest deadline, flow 1 in the two-flow case, is never reprofiled.



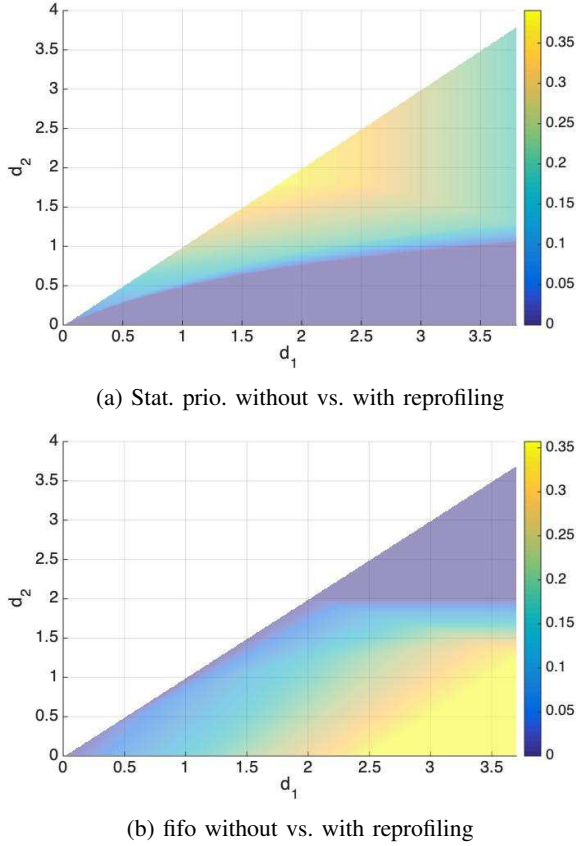


Fig. 3: Relative bandwidth increases for  $(r_1, b_1) = (4, 10)$  and  $(r_2, b_2) = (10, 18)$  as a function of  $d_1$  and  $d_2 < d_1$ . The figure is in the form of a heat-map. Darker colors (purple) correspond to smaller increases than lighter ones (yellow).

the required bandwidth only when  $d_2$ , the smaller deadline, is small. This is again intuitive as a large  $d_2$  means that the deadline can be met even without reprofiling flow 1<sup>7</sup>. Fig. 3b presents the relative gain in bandwidth for again the same 2-flow combination. As in the static priority case, the figure shows that for a fifo scheduler the benefits of reprofiling can reach about 40% in the example under consideration.

The next section explores scenarios involving combinations of multiple flow profiles. Based on those results it appears that, unsurprisingly, a fifo scheduler stands to generally benefit more from reprofiling than a static priority one.

### B. Relative Performance – Multiple Flows

In this section, we extend the investigation to configurations with more than two flows, using both synthetic flow profiles and profiles derived from datacenter traffic traces. The evaluation relies on generating a set of flow profiles, *i.e.*, token buckets plus deadlines, and for each combination compute the bandwidth required to meet the deadlines using the results derived in the paper. The main difference with the 2-flow configurations of Section VI-A is that, as described next, we

now consider a wider range of token bucket parameters with different possible combinations of deadlines. Additionally, unlike the 2-flow configurations for which the amount of bandwidth required could be obtained from explicit expressions, *i.e.*, Eq. (11), Eq. (13), and Eq. (15), computing the required bandwidth now typically involves numerical procedures, as documented in the propositions derived in the paper.

1) *Synthetic Flow Profiles*: We assign flows to *ten* different deadline classes with a dynamic range of 10, *i.e.*, with minimum and maximum deadlines of 0.1 and 1, respectively, and consider different spreads in that range for the 10 deadline classes. Specifically, we select three different possible types of spreads for deadline classes, namely,

*Even* deadline spread:

$$1) \mathbf{d}_{11} = (1, 0.9, 0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2, 0.1);$$

*Bi-modal* deadline spread:

$$2) \mathbf{d}_{21} = (1, 0.95, 0.9, 0.85, 0.8 || 0.3, 0.25, 0.2, 0.15, 0.1),$$

$$3) \mathbf{d}_{22} = (1, 0.96, 0.93, 0.9, 0.86, 0.83, 0.8 || 0.2, 0.15, 0.1),$$

$$4) \mathbf{d}_{23} = (1, 0.95, 0.9 || 0.3, 0.26, 0.23, 0.2, 0.16, 0.13, 0.1);$$

*Tri-modal* deadline spread:

$$5) \mathbf{d}_{31} = (1, 0.95, 0.9 || 0.6, 0.55, 0.5, 0.45 || 0.2, 0.15, 0.1),$$

$$6) \mathbf{d}_{32} = (1 || 0.68, 0.65, 0.62, 0.6, 0.57, 0.55, 0.53, 0.5 || 0.1),$$

$$7) \mathbf{d}_{33} = (1 || 0.6 || 0.28, 0.25, 0.23, 0.2, 0.17, 0.15, 0.12, 0.1),$$

$$8) \mathbf{d}_{34} = (1, 0.97, 0.95, 0.93, 0.9, 0.88, 0.85, 0.82 || 0.6 || 0.1).$$

Those three types of spreads translate into different groupings of deadlines, which affect the relative numbers of deadlines in close proximity to each others.

Each of the above eight groupings is used across 1,000 experiments, where an experiment consists of randomly selecting a “flow’s” traffic profile for each of the ten deadline classes. Note that what we denote by a flow, in practice maps to the aggregate of all individual flows assigned to the corresponding deadline class (individual flow profiles add up). Flow profiles are generated by independently drawing ten (aggregate) flow burst sizes  $b_1$  to  $b_{10}$  from  $U(1, 10)$ , and ten (aggregate) rates  $r_1$  to  $r_{10}$  from  $U(0, r_{\max})$ . The upper bound  $r_{\max}$  corresponds to a rate value beyond which a fifo scheduler always performs as well as the optimal solution even without reprofiling<sup>8</sup>.

The primary purpose of those synthetic experiments is to allow a systematic exploration of the performance of the different schemes across a broad range of configurations. The results can then be used to assess the expected performance of each scheme for individual configurations of practical interest.

The results of the experiments are summarized in Table I, which gives the mean, standard deviation, and the mean’s 95% confidence interval for the *relative savings* in link bandwidth, first for edf over static priority with reprofiling, followed by edf over fifo with reprofiling, and then static priority over fifo both with reprofiling. As mentioned, bandwidth values are computed numerically for each configuration using results from the previously derived propositions.

The first conclusion one can draw from Table I is that while an edf scheduler affords some benefits, they are on average smaller than the maximum values of Section VI-A. Average improvements over static priority with reprofiling hover around

<sup>7</sup>Note that under static priority the smaller deadline flow is reprofiled, while in the fifo case it is the larger deadline flow that is reprofiled to minimize its impact on the one with the tighter deadline.

<sup>8</sup>This happens when the sum of the rates is large enough to alone clear the aggregate burst before the smallest deadline, *i.e.*,  $10r_{\max} = \frac{\sum_{i=1}^{10} b_i}{0.1}$ .

TABLE I: Bandwidth savings from scheduler choice. Synthetic flow profiles.

Comparisons	Scenario	Mean	Std. Dev.	95% Conf.
edf vs. static+reprofiling $R^*$ vs. $\hat{R}_R^*$ $\left(\frac{\hat{R}_R^* - R^*}{\hat{R}_R^*}\right)$	$d_{11}$	1.2%	2.3%	[1.02%, 1.31%]
	$d_{21}$	1.5%	2.7%	[1.35%, 1.69%]
	$d_{22}$	1.1%	2.7%	[1.01%, 1.28%]
	$d_{23}$	2.9%	4.2%	[2.59%, 3.12%]
	$d_{31}$	1.4%	2.5%	[1.2%, 1.51%]
	$d_{32}$	1.0%	2.1%	[0.84%, 1.1%]
	$d_{33}$	6.2%	6.5%	[5.76%, 6.58%]
	$d_{34}$	0.7%	1.7%	[0.6%, 0.81%]
edf vs. fifo+reprofiling $R^*$ vs. $\hat{R}_R^*$ $\left(\frac{\hat{R}_R^* - R^*}{\hat{R}_R^*}\right)$	$d_{11}$	1.7%	6.5%	[1.13%, 2.11%]
	$d_{21}$	3.2%	8.7%	[2.68%, 3.76%]
	$d_{22}$	1.7%	6.2%	[1.26%, 2.03%]
	$d_{23}$	8.0%	12.8%	[7.24%, 8.82%]
	$d_{31}$	2.5%	7.8%	[2.06%, 3.03%]
	$d_{32}$	0.8%	4.6%	[0.54%, 1.11%]
	$d_{33}$	12.0%	14.1%	[11.15%, 12.9%]
	$d_{34}$	0.4%	3.2%	[0.2%, 0.6%]
static vs. fifo both w/ reprofiling $\hat{R}_R^*$ vs. $\hat{R}_R^*$ $\left(\frac{\hat{R}_R^* - \hat{R}_R^*}{\hat{R}_R^*}\right)$	$d_{11}$	0.6%	6.5%	[0.16%, 0.95%]
	$d_{21}$	1.8%	8.3%	[1.26%, 2.28%]
	$d_{22}$	0.5%	6.1%	[0.12%, 0.88%]
	$d_{23}$	5.5%	11.3%	[4.84%, 6.24%]
	$d_{31}$	1.2%	7.5%	[0.76%, 1.69%]
	$d_{32}$	-0.2%	4.5%	[-0.43%, 0.13%]
	$d_{33}$	6.6%	11.2%	[5.92%, 7.3%]
	$d_{34}$	-0.3%	3.3%	[-0.53%, -0.12%]

1% and did not exceed about 6%. Improvements are a little higher when considering fifo with reprofiling, where they reach 12%, but those values are still significantly less than the worst case scenarios of Section VI-A.

Table I also reveals that, somewhat surprisingly, static priority and fifo perform similarly when both are afforded the benefit of reprofiling (the largest difference observed in the experiments is 5.5%). Static priority has an edge on average even if, as discussed in Section VI-A, a few scenarios exist where a fifo scheduler outperforms static priority when both are combined with reprofiling, e.g.,  $d_{32}$  and  $d_{34}$ . Recall that this is because we strictly map smaller deadlines to higher priority. The differences are, however, small, i.e., 0.2% and 0.3%, respectively, for the two scenarios where fifo outperforms static priority on average.

TABLE II: Benefits of reprofiling for static priority & fifo. Synthetic flow profiles.

Comparisons	Scenario	Mean	Std. Dev.	95% Conf.
static w/ & w/o reprofiling $\hat{R}_R^*$ vs. $\hat{R}_R^*$ $\left(\frac{\hat{R}_R^* - \hat{R}_R^*}{\hat{R}_R^*}\right)$	$d_1$	8.43%	4.50%	[8.15%, 8.71%]
	$d_{21}$	8.11%	4.19%	[7.85%, 8.37%]
	$d_{22}$	8.42%	4.52%	[8.14%, 8.71%]
	$d_{23}$	9.38%	4.80%	[9.08%, 9.67%]
	$d_{31}$	8.24%	4.33%	[7.97%, 8.51%]
	$d_{32}$	9.49%	5.07%	[9.18%, 9.81%]
	$d_{33}$	15.97%	4.78%	[15.67%, 16.27%]
	$d_{34}$	8.83%	4.94%	[8.53%, 9.14%]
fifo w/ & w/o reprofiling $\hat{R}_R^*$ vs. $\hat{R}_R^*$ $\left(\frac{\hat{R}_R^* - \hat{R}_R^*}{\hat{R}_R^*}\right)$	$d_1$	49.52%	8.17%	[49.01%, 50.03%]
	$d_{21}$	48.71%	7.62%	[48.24%, 49.18%]
	$d_{22}$	49.53%	8.27%	[49.02%, 50.05%]
	$d_{23}$	45.78%	6.52%	[45.37%, 46.18%]
	$d_{31}$	49.08%	7.88%	[48.59%, 49.57%]
	$d_{32}$	49.95%	8.59%	[49.42%, 50.49%]
	$d_{33}$	42.47%	6.19%	[42.08%, 42.85%]
	$d_{34}$	50.13%	8.84%	[49.59%, 50.68%]

Towards gaining a better understanding of reprofiling and the extent to which it is behind the somewhat unexpected

good performance of fifo, Table II reports its impact for both static priority and fifo. As Table I, it gives the mean, standard deviation, and the mean's 95% confidence interval, but now of the relative gains in bandwidth that reprofiling affords over no reprofiling for both fifo and static priority schedulers.

The data from Table II highlights that while both static priority and fifo benefit from reprofiling, the magnitude of the improvements is significantly higher for fifo. Specifically, improvements from reprofiling are systematically above 40% and often close to 50% for fifo, while they exceed 10% only once for static priority (at 15% for scenario  $d_{33}$ ) and are typically around 8%. As alluded to earlier, this is not surprising given that static priority offers some ability to discriminate flows based on their deadlines, while fifo does not.

2) *Application Derived Flow Profiles*: The benefits of synthetic flow profiles in allowing a systematic investigation notwithstanding, it is also of interest to target configurations more directly representative of traffic mixes as they arise in practice. To that end, we rely on a methodology similar to that used in [17, Section VIII-B2], and construct a set of flow profiles derived from traffic data reported in [24].

Specifically, [24] investigates the traffic flowing through the network of one of Facebook's large datacenter, and reports, among other things, the distribution of flow sizes and durations (Figs. 6 & 7 of [24]) for three representative applications: Web (W), Cache read and replacement (C), and Hadoop (H). We rely on these data to generate sample traffic profiles ( $r, b$ ) for flows from those three applications as follows:

- 1) For a given application, we generate flow size+duration tuples by sampling the corresponding distributions assuming they are perfectly positively correlated. In other words, we assume that larger flows last longer.
- 2) A flow's token rate  $r$  is then obtained by dividing the flow size by its duration. Fig. 4 shows the resulting cumulative distributions of flow rates for all three applications.
- 3) Generating token bucket sizes  $b$  involves an additional step and associated assumption:
  - a) The smallest flow sizes from Fig. 6 of [24] are assumed representative of a single transmission burst. This yields burst sizes  $S_W = 0.15$ Kbytes,  $S_C = 0.4$ Kbytes, and  $S_H = 0.3$ Kbytes for our three sample applications.
  - b) As bucket sizes are typically chosen to accommodate consecutive bursts, we leverage the claim in [24] that all three applications are "internally bursty" with Cache significantly burstier than Hadoop, and Web in between, to randomly select bucket sizes in  $[0, 20S_C]$ ,  $[0, 10S_W]$  and  $[0, 2S_H]$ , respectively. We note that these values yield relatively small buckets, and, therefore, maximum burst sizes.

The resulting profiles have relatively low rates and burstiness, at least when it comes to individual flows, with Hadoop's profile typical of bandwidth hungry applications, and Web and Cache representative of more interactive applications. This maps to the types of services that [24] mentions as relying on those three applications, i.e., Web search, user data query, and offline analysis (e.g., data mining). As a result, we assign deadlines to each application that broadly reflect those services, with three deadline classes set to 10ms, 50ms, and

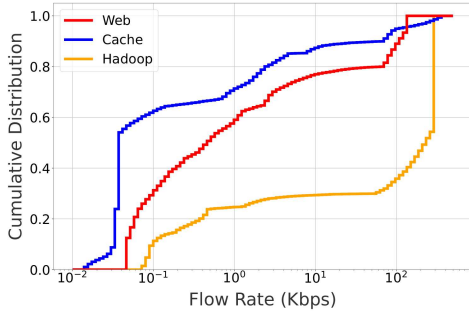


Fig. 4: CDF of flow rates for Web, Cache, & Hadoop applications from [24] assuming correlated flow durations and sizes.

200ms for Web, Cache, and Hadoop respectively.

In evaluating performance, we consider four “traffic mixes” that differ in their relative proportion of flows from each application. Specifically, the corresponding four scenarios sample our W, C, H applications in the proportions: 1:1:1, 3:9:1, 9:3:1, and 9:9:1, respectively. For each scenario, we randomly sample 100 flow profiles in those proportions. The 100 flows are then grouped according to their deadline class, which yields a set of three aggregates to schedule on the shared link according to their deadlines. This procedure is repeated 1000 times, and the relative link bandwidth requirements across schedulers and reprofiling options are given in Table III.

TABLE III: Bandwidth savings & benefits of reprofiling. Application-derived flow profiles.

Comparisons	Scenario	Mean	Std. Dev.	95% Conf.
$R^*$ vs. $\tilde{R}_R^*$ $\left(\frac{\tilde{R}_R^* - R^*}{\tilde{R}_R^*}\right)$	1:1:1	0.0%	0.0%	[0%, 0%]
	3:9:1	0.0%	0.0%	[0%, 0%]
	9:3:1	0.0%	0.0%	[0%, 0%]
	9:9:1	0.0%	0.0%	[0%, 0%]
$R^*$ vs. $\hat{R}_R^*$ $\left(\frac{\hat{R}_R^* - R^*}{\hat{R}_R^*}\right)$	1:1:1	41.82%	17.04%	[40.76%, 42.88%]
	3:9:1	72.36%	3.43%	[72.15%, 72.58%]
	9:3:1	56.84%	8.43%	[56.32%, 57.37%]
	9:9:1	69.70%	4.10%	[69.45%, 69.96%]
$\tilde{R}_R^*$ vs. $\hat{R}_R^*$	See above			
$\tilde{R}_R^*$ vs. $\tilde{R}^*$ $\left(\frac{\tilde{R}^* - \tilde{R}_R^*}{\tilde{R}^*}\right)$	1:1:1	0.01%	0.22%	[0.00%, 0.03%]
	3:9:1	1.74%	0.73%	[1.70%, 1.79%]
	9:3:1	1.02%	2.34%	[0.87%, 1.16%]
	9:9:1	4.06%	1.19%	[3.98%, 4.13%]
$\hat{R}_R^*$ vs. $\hat{R}^*$ $\left(\frac{\hat{R}^* - \hat{R}_R^*}{\hat{R}^*}\right)$	1:1:1	22.74%	9.17%	[22.17%, 23.31%]
	3:9:1	21.61%	8.48%	[21.09%, 22.14%]
	9:3:1	14.93%	8.93%	[14.37%, 15.48%]
	9:9:1	19.55%	8.62%	[19.02%, 20.08%]

The first conclusion from Table III is that static priority with reprofiling performs just as well as edf for all four scenarios (top four rows). This is not surprising. The large gaps between the deadlines of the three classes of applications ensure that once high-priority bursts are cleared, the residual bandwidth is sufficient to transmit lower priority bursts before their deadline. As with synthetic profiles, reprofiling is instrumental in realizing this outcome, even if the wide gaps between deadlines together with the limited burstiness of the applications produce a smaller gain ( $\tilde{R}_R^*$  vs.  $\hat{R}^*$ ).

The benefits of reprofiling are again more apparent with fifo, even if it significantly under-performs both edf and static priority. The lack of discrimination across flows that

fifo suffers from is exacerbated by the large gaps between deadlines, and compensating for it calls for an average of about 50% more bandwidth across all four scenarios. However, without reprofiling this bandwidth increase ( $\hat{R}_R^*$  vs.  $\hat{R}^*$ ) is around 20% larger. This again demonstrates the extent to which reprofiling can help simpler schedulers.

### C. Summary Discussion

Several common themes emerge between the evaluations of Sections VI-B1 and VI-B2. The first is that reprofiling can help a static priority scheduler perform nearly as well as an edf scheduler. Second, while its inability to discriminate between flows puts fifo at a clear disadvantage, reprofiling is again capable of partially mitigating its handicap. Finally, while with static priority the benefits of reprofiling are realized by reprofiling high-priority flows to limit their impact on low-priority ones, the opposite holds for fifo.

The differences between the results of Sections VI-B1 and VI-B2 also revealed a number of intuitive findings brought about by the differences in deadline spreads in the two scenarios. In particular, the large gaps between deadlines present in Section VI-B2 make it easier for a static priority scheduler to perform nearly as well as an edf scheduler. This holds with and without reprofiling, even if reprofiling remains useful. Conversely, more closely packed deadlines offer additional opportunities for reprofiling to be useful, as closer deadlines amplify the need for fine tuning of a flow’s profile relative to its deadline and impact on other flows.

## VII. RELATED WORKS

The question of meeting deadlines for a set of rate-limited flows is one that has received much attention in the scheduling literature. It is not our intent to provide an exhaustive review of those works. Instead, we limit ourselves to highlighting works whose results are closest to ours or that offered early insight into the problem, including the benefits of adjusting flows’ profiles (reprofiling) that is one of the foci of this paper.

a) *Packet-level shaping and scheduling*: Scheduling flows with deterministic traffic profiles was investigated in [19] that considered both buffer and delay requirements. In particular, the paper established<sup>9</sup> the optimality of the edf policy in terms of maximizing the schedulable region. This is the “dual” of the bandwidth minimization problem investigated in this paper, and the result parallels that of Proposition 2. Static priority and fifo schedulers were, however, not investigated, and neither was the impact of reprofiling flows.

The aspect of minimizing the resources required to meet the latency targets of token bucket-controlled flows was explored in [26]. The paper relied on service curves with high and low rates and sought to identify the earliest possible time for switching to the lower rate. The focus was, however, on minimizing resources required by each flow individually rather than in aggregate, as in this paper. In addition, the potential impact of reprofiling flows was not addressed.

<sup>9</sup>A similar result was reported in [25].

b) *Shaping bulk data transfers*: Minimizing bandwidth (cost) through reprofiling (reshaping) flows has been investigated for bulk data transfers where transfer completion times rather than packet-level deadlines are the targets [27]–[32]. The problem stems from non-linear bandwidth costs, e.g., based on the 95<sup>th</sup> percentile, so that judiciously adjusting (shaping) the transmission rates of bulk transfers can yield significant savings. Rate shaping is, however, at a time-granularity of minutes rather than at the packet-level. The optimization frameworks of those papers are, therefore, not applicable to our problem. Their solutions, can, however, complement ours by leveraging the fluctuations in link utilization inherent in delivering hard, packet-level delay bounds, as we do.

c) *Deterministic networking*: The deterministic traffic profiles and delay bounds of the TSN and DetNet standards have also given rise to related investigations as documented in recent surveys [33], [34]. In particular, the optimization framework that underlies many of those studies have connections to the problem we address. However, like most prior similar works, traffic profiles are assumed fixed and the impact of reprofiling is not considered.

d) *Datacenter solutions*: The emergence of traffic profiles and latency targets in datacenter networks motivated [35]. It targets a multi-hop network, but calls on topological properties of typical datacenter networks to collapse its model to a single hop, thereby aligning with the scope of this paper. Similarities extend to considering a static priority scheduler, but traffic profiles differ. Rather than a token bucket, [35] relies on the notion of network “epochs” to bound packet bursts. Delay bounds are then expressed as a function of the network fan-in and a “throughput factor” that reflects the number of transmission opportunities sources can have per network epoch. Also absent from the paper are exploring bandwidth minimization and the potential benefits of reprofiling flows.

e) *Reprofiling investigations*: Meeting packet-level latency constraints with a static priority scheduler while minimizing costs through reprofiling of flows is the focus of WorkloadCompactor [15]. The reprofiling decisions of [15] are, however, focused on selecting token bucket parameters from among a family of feasible regulators<sup>10</sup> that do not introduce additional delay. In contrast, our reprofiling allows for an added delay that the scheduler must then compensate for. Exploring when and how this trade-off is of benefit is our main contribution and what makes this work *complementary* to the approach from WorkloadCompactor.

Specifically, WorkloadCompactor considers traffic/workload traces for which it seeks to first identify *feasible* token bucket parameter pairs  $\langle r, b \rangle$  that result in *zero* access delay for those traces. WorkloadCompactor’s main contribution is in realizing that multiple such  $\langle r, b \rangle$  pairs are possible (the  $r$ - $b$  curve of [15]), and that “*jointly optimizing the choice of  $\langle r, b \rangle$  rate limit parameters for each workload to better compact workloads onto servers*” can reduce the required server capacity. This is where the contribution of WorkloadCompactor ends, and where that of this paper actually starts.

More precisely, once the optimization of WorkloadCompactor completes, the set of  $\langle r, b \rangle$  pairs it produces can be used, together with the associated target latency bounds, as inputs to Proposition 8. Proposition 8 explores, for a static priority scheduler, how to best reprofile token bucket-controlled flows to meet their deadlines with the least bandwidth. This reprofiling is beyond that suggested by WorkloadCompactor<sup>11</sup>, and explores how trading-off access delay to further smooth flows can yield additional benefits. In other words, the approach proposed in the paper complements that of WorkloadCompactor, in that it can be applied to any set of token bucket profiles produced by WorkloadCompactor, and modify them to yield further reductions in system resources (bandwidth or server capacity) while still meeting latency targets.

We also note that, because WorkloadCompactor considers the problem of selecting token bucket parameters for traffic traces (to ensure zero access delay), it addresses an aspect that this paper does not consider since we assume that token bucket profiles are given. This is yet another aspect in which the two papers are complementary.

f) *Early works*: Finally, we note that exploring the trade-off between making traffic smoother and end-to-end performance is not unique to packet networks. It is present in the early “fluctuation smoothing” scheduling policies of [36] that sought to reduce processing time in manufacturing plants, and more recently in the reshaping of parallel I/O requests to improve the scalability of database systems [37].

## VIII. CONCLUSION AND FUTURE WORK

The paper investigated the question of minimizing the bandwidth needed to guarantee worst case latencies to a set of token bucket-controlled flows sharing a single link. The investigation was carried for schedulers of different complexity.

The paper first characterized the minimum required bandwidth independent of schedulers, and showed that an edf scheduler could realize all flows’ deadlines under such bandwidth. Motivated by the need for lower complexity solutions, the paper then explored simpler static priority and fifo schedulers. It derived the minimum required bandwidth for both, but more interestingly established how to optimally reprofile flows to reduce the bandwidth needed while still meeting all deadlines. The relative benefits of such an approach were illustrated numerically for a number of different flow combinations, which showed how reprofiling can enable simpler schedulers to perform nearly as well a more complex ones across a range of configurations.

The obvious direction in which to extend the paper is to a multi-hop setting. In [17], we build on the results of Proposition 1 and provide initial results for the multi-hop case under the assumption that (service curve) edf schedulers are available at each hop. Extending the investigation to static priority and fifo schedulers is under way.

Another aspect of interest with static priority schedulers is relaxing the assumption that flows with different deadlines map to distinct priority classes, and allow multiple deadlines

<sup>10</sup>Regulators above the  $r$ - $b$  curve using the terminology of [15].

<sup>11</sup>Again [15] focuses on regulators that ensure zero access delay for a given trace, while we investigate how a non-zero access delay can be of benefit.

to be assigned to the same class. Not only does it enhance scalability, but it can also improve performance<sup>12</sup>. Last but not least, extensions to statistical rather than deterministic delay guarantees are also of practical relevance.

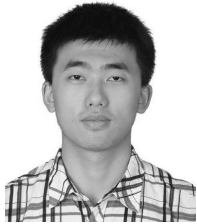
## REFERENCES

- [1] M. Ashjaei, L. L. Bello, M. Daneshmand, G. Patti, S. Saponara, and S. Mubeen, "Time-sensitive networking in automotive embedded systems: State of the art and research opportunities," *Journal of Systems Architecture*, vol. 117, p. 102137, 2021. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1383762121001028>
- [2] *Avionics Full Duplex Switched Ethernet (AFDX) Network*, Airlines Electronic Engineering Committee, Aircraft Data Network Part 7, ARINC Specification 664, Aeronautical Radio, Annapolis, MD, USA, 2002.
- [3] C. Zunino, A. Valenzano, R. Obermaier, and S. Petersen, "Factory communications at the dawn of the fourth industrial revolution," *Computer Standards & Interfaces*, vol. 71, p. 103433, 2020. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0920548919300868>
- [4] T. Docquier, Y.-Q. Song, V. Chevrier, L. Pontnau, and A. Ahmed-Nacer, "Performance evaluation methodologies for smart grid substation communication networks: A survey," *Computer Communications*, vol. 198, pp. 228–246, 2023. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0140366422004285>
- [5] (2022) AWS global network. [Online]. Available: <https://aws.amazon.com/about-aws/global-infrastructure/global-network/>
- [6] (2022) Google cloud networking overview. [Online]. Available: <https://cloud.google.com/blog/topics/developers-practitioners/google-cloud-networking-overview>
- [7] (2021, January) Microsoft global network. [Online]. Available: <https://docs.microsoft.com/en-us/azure/networking/microsoft-global-network>
- [8] J. Farkas, L. L. Bello, and C. Gunther, "Time-sensitive networking standards," *IEEE Communications Standards Magazine*, vol. 2, no. 2, pp. 20–21, June 2018.
- [9] G. Parsons, "The rise of time-sensitive networking (TSN) in automobiles, industrial automation, and aviation," In *Compliance - Electronic Design, Testing & Standards*, January 2022, <https://incompliancemag.com/article/the-rise-of-time-sensitive-networking-tsn-in-automobiles-industrial-automation-and-aviation/>.
- [10] Y. Seol, D. Hyeon, J. Min, M. Kim, and J. Paek, "Timely survey of time-sensitive networking: Past and future directions," *IEEE Access*, vol. 9, pp. 142 506–142 527, 2021.
- [11] N. Finn, P. Thubert, B. Varga, and J. Farkas, "Deterministic Networking Architecture," RFC 8655, October 2019. [Online]. Available: <https://www.rfc-editor.org/info/rfc8655>
- [12] C. Jones, J. Wilkes, N. Murphy, and C. Smith, *Site Reliability Engineering: How Google Runs Production Systems*, B. Beyer, C. Jones, J. Petoff, and N. Murphy, Eds. O'Reilly Media, 2016, accessible at <https://sre.google/sre-book/service-level-objectives/>.
- [13] (2022) Enterprises are willing to pay for an SLA with guaranteed latency, says new 5G MEC report. [Online]. Available: <https://www.vanillaplus.com/2022/06/08/70101-enterprises-are-willing-to-pay-for-an-sla-with-guaranteed-latency-says-new-5g-mec-report/>
- [14] J.-Y. Le Boudec, "A theory of traffic regulators for deterministic networks with application to interleaved regulators," *IEEE/ACM Transactions on Networking*, vol. 26, no. 6, p. 2721–2733, December 2018. [Online]. Available: <https://doi.org/10.1109/TNET.2018.2875191>
- [15] T. Zhu, M. A. Kozuch, and M. Harchol-Balter, *WorkloadCompactor: Reducing Datacenter Cost While Providing Tail Latency SLO Guarantees*. Santa Clara, CA: Association for Computing Machinery, September 2017, p. 598–610. [Online]. Available: <https://doi.org/10.1145/3127479.3132245>
- [16] J. Song, J. Qiu, R. Guérin, and H. Sariowan, "On the benefits of traffic "reprofiling" – The single hop case," 2024. [Online]. Available: <https://arxiv.org/abs/2104.02222>
- [17] J. Qiu, J. Song, R. Guérin, and H. Sariowan, "On the benefits of traffic "reprofiling." The multiple hops case – Part I," 2023, under submission.
- [18] J.-Y. Le Boudec and P. Thiran, *Network Calculus: A Theory of Deterministic Queuing Systems for the Internet*. Springer, 2001.
- [19] L. Georgiadis, R. Guérin, and A. Parekh, "Optimal multiplexing on a single link: delay and buffer requirements," *IEEE Transactions on Information Theory*, vol. 43, no. 5, pp. 1518–1535, September 1997.
- [20] A. K. Parekh and R. G. Gallager, "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, June 1993.
- [21] A. Sivaraman, S. Subramanian, M. Alizadeh, S. Chole, S.-T. Chuang, A. Agrawal, H. Balakrishnan, T. Edsall, S. Katti, and N. McKeown, "Programmable packet scheduling at line rate," in *Proc. ACM SIGCOMM*, Florianopolis, Brazil, August 2016.
- [22] A. Saeed, Y. Zhao, N. Dukkupati, M. Ammar, E. Zegura, K. Harras, and A. Vahdat, "Eiffel: Efficient and flexible software packet scheduling," in *Proc. USENIX NSDI*, Boston, MA, 2019.
- [23] N. K. Sharma, C. Zhao, M. Liu, P. G. Kannan, C. Kim, A. Krishnamurthy, and A. Sivaraman, "Programmable calendar queues for high-speed packet scheduling," in *Proc. USENIX NSDI*, Santa Clara, CA, February 2020.
- [24] A. Roy, H. Zeng, J. Bagga, G. Porter, and A. C. Snoeren, "Inside the social network's (datacenter) network," in *Proc. ACM SIGCOMM*, London, United Kingdom, August 2015. [Online]. Available: <https://doi.org/10.1145/2785956.2787472>
- [25] J. Liebeherr, D. E. Wrege, and D. Ferrari, "Exact admission control for networks with a bounded delay service," *IEEE/ACM Transactions on Networking*, vol. 4, no. 6, pp. 885–901, 1996.
- [26] J. B. Schmitt, "On the allocation of network service curves for bandwidth/delay-decoupled scheduling disciplines," in *Proc. IEEE GLOBECOM*, vol. 2, Taipei, Taiwan, November 2002.
- [27] T. Nandagopal and K. P. Puttaswamy, "Lowering inter-datacenter bandwidth costs via bulk data scheduling," in *Proc. IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing*, Ottawa, ON, Canada, May 2012.
- [28] H. Zhang, K. Chen, W. Bai, D. Han, C. Tian, H. Wang, H. Guan, and M. Zhang, "Guaranteeing deadlines for inter-datacenter transfers," in *Proc. European Conference on Computer Systems (EuroSys)*, Bordeaux, France, 2015. [Online]. Available: <https://doi.org/10.1145/2741948.2741957>
- [29] V. Jalaparti, I. Bliznets, S. Kandula, B. Lucier, and I. Menache, "Dynamic pricing and traffic engineering for timely inter-datacenter transfers," in *Proc. ACM SIGCOMM*, Florianopolis, Brazil, August 2016. [Online]. Available: <https://doi.org/10.1145/2934872.2934893>
- [30] W. Li, X. Zhou, K. Li, H. Qi, and D. Guo, "TrafficShaper: Shaping inter-datacenter traffic to reduce the transmission cost," *IEEE/ACM Transactions on Networking*, vol. 26, no. 3, pp. 1193–1206, June 2018.
- [31] Z. Yang, Y. Cui, X. Wang, Y. Liu, M. Li, S. Xiao, and C. Li, "Cost-efficient scheduling of bulk transfers in inter-datacenter WANs," *IEEE/ACM Transactions on Networking*, vol. 27, no. 5, pp. 1973–1986, 2019.
- [32] R. Singh, S. Agarwal, M. Calder, and P. Bahl, "Cost-effective cloud edge traffic engineering with cascara," in *Proc. USENIX NSDI*, virtual, 2021. [Online]. Available: <https://www.usenix.org/conference/nsdi21/presentation/singh>
- [33] T. Stüber, L. Osswald, S. Lindner, and M. Menth, "A survey of scheduling in time-sensitive networking (TSN)," 2022. [Online]. Available: <https://arxiv.org/abs/2211.10954>
- [34] J. Walrand, "A concise tutorial on traffic shaping and scheduling in time-sensitive networks," *IEEE Communications Surveys & Tutorials*, pp. 1–1, May 2023.
- [35] M. P. Grosvenor, M. Schwarzkopf, I. Gog, R. N. M. Watson, A. W. Moore, S. Hand, and J. Crowcroft, "Queues Don't matter when you can JUMP them!" in *Proc. USENIX NSDI*, Oakland, CA, May 2015. [Online]. Available: <https://www.usenix.org/conference/nsdi15/technical-sessions/presentation/grosvenor>
- [36] S. Lu, D. Ramaswamy, and P. Kumar, "Efficient scheduling policies to reduce mean and variance of cycle-time in semiconductor manufacturing plants," *IEEE Transactions on Semiconductor Manufacturing*, vol. 7, no. 3, pp. 374–388, 1994.
- [37] N. Li, H. Jiang, H. Che, Z. Wang, and M. Q. Nguyen, "Improving scalability of database systems by reshaping user parallel I/O," in *Proc. EuroSys*, Rennes, France, April 2022. [Online]. Available: <https://doi.org/10.1145/3492321.3519570>

<sup>12</sup>See Appendix E of [16] for the simple case of a two-flow configuration.



**Jiayi Song** received her Ph.D. in computer science from Washington University in St. Louis. Prior to her doctoral studies, she obtained her bachelor's degree from Wuhan University, majoring in both applied mathematics and economics. Currently, she is with Bytedance Inc.'s Technical Infrastructures team.



**Jiaming Qiu** received his B.S. in 2020 and is currently a Ph.D candidate, both in Computer Science from Washington University in St. Louis. His research interests include network performance evaluation and optimization, edge computing, and machine learning.



**Roch Guérin** (F IEEE '01 / F ACM '06) received his Ph.D. from Caltech, worked at the IBM T. J. Watson Research Center and then the University of Pennsylvania. He is currently with the Department of Comp. Sci. & Eng. of Washington University in St. Louis, as the Harold B. and Adelaide G. Welge Professor and department chair. Dr. Guérin was the Editor-in-Chief for the IEEE/ACM Transactions on Networking from 2009 till 2013 and served as chair of ACM SIGCOMM. In 2010 he received the IEEE INFOCOM Achievement Award for "Pioneering Contributions to the Theory and Practice of QoS in Networks".



**Henry Sariowan** (SM IEEE, '03) received his Ph.D. from the University of California, San Diego, M.S. from Columbia University, New York, and B.S. from Sepuluh Nopember Institute of Technology, Surabaya, Indonesia, all in electrical engineering. He is currently with Google's Global Networking team, part of Google Cloud. Dr. Sariowan previously worked with several technology companies in San Diego, CA and published papers in the areas of network quality of service and video transmission.