



Universal Packet Scheduling

Radhika Mittal, Rachit Agarwal, and Sylvia Ratnasamy, *University of California, Berkeley*;
Scott Shenker, *University of California, Berkeley, and International Computer Science Institute*

<https://www.usenix.org/conference/nsdi16/technical-sessions/presentation/mittal>

This paper is included in the Proceedings of the
13th USENIX Symposium on Networked Systems
Design and Implementation (NSDI '16).

March 16–18, 2016 • Santa Clara, CA, USA

ISBN 978-1-931971-29-4

Open access to the Proceedings of the
13th USENIX Symposium on
Networked Systems Design and
Implementation (NSDI '16)
is sponsored by USENIX.

Universal Packet Scheduling

Radhika Mittal[†]

Rachit Agarwal[†]
[†]UC Berkeley

Sylvia Ratnasamy[†]
[‡]ICSI

Scott Shenker^{†‡}

Abstract

In this paper we address a seemingly simple question: *Is there a universal packet scheduling algorithm?* More precisely, we analyze (both theoretically and empirically) whether there is a single packet scheduling algorithm that, at a network-wide level, can perfectly match the results of *any* given scheduling algorithm. We find that in general the answer is “no”. However, we show theoretically that the classical Least Slack Time First (LSTF) scheduling algorithm comes closest to being universal and demonstrate empirically that LSTF can closely replay a wide range of scheduling algorithms. We then evaluate whether LSTF can be used *in practice* to meet various network-wide objectives by looking at popular performance metrics (such as average FCT, tail packet delays, and fairness); we find that LSTF performs comparable to the state-of-the-art for each of them. We also discuss how LSTF can be used in conjunction with active queue management schemes (such as CoDel and ECN) without changing the core of the network.

1 Introduction

There is a large and active research literature on novel packet scheduling algorithms, from simple schemes such as priority scheduling [38], to more complicated mechanisms for achieving fairness [20, 34, 39], to schemes that help reduce tail latency [19] or flow completion time [10], and this short list barely scratches the surface of past and current work. In this paper we do not add to this impressive collection of algorithms, but instead ask if there is a single *universal* packet scheduling algorithm that could obviate the need for new ones. In this context, we consider a packet scheduling algorithm to be both how packets are served inside the network (based on their arrival times and their packet headers) and how packet header fields are initialized and updated; this definition includes all the classical scheduling algorithms (FIFO, LIFO, priority, round-robin) as well as algorithms that incorporate dynamic packet state [19, 44, 45].

We can define a universal packet scheduling algorithm (hereafter UPS) in two ways, depending on our viewpoint on the problem. From a theoretical perspective, we call a packet scheduling algorithm *universal* if it can replay any *schedule* (the set of times at which packets arrive to and exit from the network) produced by any other scheduling algorithm. This is not of practical interest, since such schedules are not typically known in advance, but it offers

a theoretically rigorous definition of universality that (as we shall see) helps illuminate its fundamental limits (i.e., which scheduling algorithms have the flexibility to serve as a UPS, and why).

From a more practical perspective, we say a packet scheduling algorithm is universal if it can achieve different desired performance objectives (such as fairness, reducing tail latencies and minimizing flow completion times). In particular, we require that the UPS should match the performance of the best known scheduling algorithm for a given performance objective.¹

The notion of universality for packet scheduling might seem esoteric, but we think it helps clarify some basic questions. If there exists no UPS then we should *expect* to design new scheduling algorithms as performance objectives evolve. Moreover, this would make a strong argument for switches being equipped with programmable packet schedulers so that such algorithms could be more easily deployed (as argued in [42]; in fact, it was the eloquent argument in this paper that caused us to initially ask the question about universality).

However, if there is indeed a UPS, then it changes the lens through which we view the design and evaluation of scheduling algorithms: e.g., rather than asking whether a new scheduling algorithm meets a performance objective, we should ask whether it is easier/cheaper to implement/configure than the UPS (which could also meet that performance objective). Taken to the extreme, one might even argue that the existence of a (practical) UPS greatly diminishes the need for programmable *scheduling* hardware.² Thus, while the rest of the paper occasionally descends into scheduling minutiae, the question we are asking has important practical (and intriguing theoretical) implications.

This paper starts from the theoretical perspective, defining a formal model of packet scheduling and our

¹For this definition of universality, we allow the header initialization to depend on the objective being optimized. That is, while the basic scheduling operations must remain constant, the header initialization can depend on whether you are seeking fairness or minimal flow completion time.

²Note that the case for programmable hardware as made in recent work on P4 and the RMT switch [14, 15] remains: these systems target programmability in header parsing and in how a packet’s processing pipeline is defined (i.e., how forwarding ‘actions’ are applied to a packet). The P4 language does not currently offer primitives for scheduling and, perhaps more importantly, the RMT switch does not implement a programmable packet scheduler; we hope our results can inform the discussion on whether and how P4/RMT might be extended to support programmable scheduling.

notion of replayability in §2. We first prove that there is no UPS, but then show that Least Slack Time First (LSTF) [28] comes as close as any scheduling algorithm to achieving universality. We also demonstrate empirically (via simulation) that LSTF can closely approximate the schedules of many scheduling algorithms. Thus, while not a perfect UPS in terms of replayability, LSTF comes very close to functioning as one.

We then take a more practical perspective in §3, showing (via simulation) that LSTF is comparable to the state of the art in achieving various objectives relevant to an application’s performance. We investigate in detail LSTF’s ability to minimize average flow completion times, minimize tail latencies, and achieve per-flow fairness. We also consider how LSTF can be used in multitenant situations to achieve multiple objectives simultaneously, while highlighting some of its key limitations.

In §4, we look at how network feedback for active queue management (AQM) can be incorporated using LSTF. Rather than augmenting the basic LSTF logic (which is restricted to packet scheduling) with a queue management algorithm, we show that LSTF can, instead, be used to implement AQM at the edge of the network. This novel approach to AQM is a contribution in itself, as it allows the algorithm to be upgraded without changing internal routers.

We then discuss the feasibility of implementing LSTF (§5) and provide an overview of related work (§6) before concluding with a discussion of open questions in §7.

2 Theory: Replaying Schedules

This section delves into the theoretical viewpoint of a UPS, in terms of its ability to *replay* a given schedule.

2.1 Definitions and Overview

Network Model: We consider a network of store-and-forward output-queued routers connected by links. The input load to the network is a fixed set of packets $\{p \in P\}$, their arrival times $i(p)$ (*i.e.*, when they reach the ingress router), and the path $path(p)$ each packet takes from its ingress to its egress router. We assume no packet drops, so all packets eventually exit. Every router executes a non-preemptive scheduling algorithm which need not be work-conserving or deterministic and may even involve oracles that know about future packet arrivals. Different routers in the network may use different scheduling logic. For each incoming load $\{(p, i(p), path(p))\}$, a collection of scheduling algorithms $\{A_\alpha\}$ (router α implements algorithm A_α) will produce a set of packet output times $\{o(p)\}$ (the time a packet p exits the network). We call the set $\{(path(p), i(p), o(p))\}$ a *schedule*.

Replaying a Schedule: Applying a different collection of scheduling algorithms $\{A'_\alpha\}$ to the same set of packets $\{(p, i(p), path(p))\}$ (with the packets taking the same path in the replay as in the original schedule), produces a new

set of output times $\{o'(p)\}$. We say that $\{A'_\alpha\}$ *replays* $\{A_\alpha\}$ on this input if and only if $\forall p \in P, o'(p) \leq o(p)$.³

Universal Packet Scheduling Algorithm: We say a schedule $\{(path(p), i(p), o(p))\}$ is *viable* if there is at least one collection of scheduling algorithms that produces that schedule. We say that a scheduling algorithm is *universal* if it can replay *all* viable schedules. While we allowed significant generality in defining the scheduling algorithms that a UPS seeks to replay (demanding only that they be non-preemptive), we insist that the UPS itself obey several practical constraints (although we allow it to be preemptive for theoretical analysis, but then quantitatively analyze the non-preemptive version in §2.3).⁴ The three practical constraints we impose on a UPS are:

(1) *Uniformity and Determinism:* A UPS must use the same deterministic scheduling logic at every router.

(2) *Limited state used in scheduling decisions:* We restrict a UPS to using only (i) packet headers, and (ii) static information about the network topology, link bandwidths, and propagation delays. It cannot rely on oracles or other external information. However, it can modify the header of a packet before forwarding it (resulting in *dynamic packet state* [45]).

(3) *Limited state used in header initialization:* We assume that the header for a packet p is initialized at its ingress node. The additional information available to the ingress for this initialization is limited to: (i) $o(p)$ from the original schedule⁵ and (ii) $path(p)$. Later, we extend the kinds of information the header initialization process can use, and find that this is a key determinant in whether one can find a UPS.

We make three observations about the above model. First, our model assumes greater capability at the edge than in the core, in keeping with common assumptions that the network edge is capable of greater processing complexity, exploited by many architectural proposals [16, 36, 44]. Second, when initializing a packet p ’s header, a UPS can only use the input time, output time and the path information for p itself, and must be *oblivious* [24] to the corresponding attributes for *other* packets in the network. Finally, the key source of impracticality in our model is the assumption that the output times $o(p)$ are known at the ingress. However,

³We allow the inequality because, if $o'(p) < o(p)$, one can delay the packet upon arrival at the egress node to ensure $o'(p) = o(p)$.

⁴The issue of preemption is somewhat complicated. Allowing the original scheduling algorithms to be preemptive allows packets to be fragmented, which then makes replay extremely difficult even in simple networks (with store-and-forward routers). However, disallowing preemption in the candidate UPS overly limits the flexibility and would again make replay impossible even in simple networks. Thus, we take the seemingly hypocritical but only theoretically tractable approach and disallow preemption in the original scheduling algorithms but allow preemption in the candidate UPS. In practice, when we care only about approximately replaying schedules, the distinction is of less importance, and we simulate LSTF in the non-preemptive form.

⁵Note that this ingress router can directly observe $i(p)$ as the time the packet arrives.

a different interpretation of $o(p)$ suggests a more practical application of replayability (and thus our results): *if we assign $o(p)$ as the “desired” output time for each packet in the network, then the existence of a UPS tells us that if these goals are viable then the UPS will be able to meet them.*

2.2 Theoretical Results

For brevity, in this section we only summarize our key theoretical results. The detailed proofs are in Appendix A.

Existence of a UPS under omniscient initialization: Suppose we give the header-initialization process extensive information in the form of times $o(p, \alpha)$ which represent when p was scheduled by router α in the original schedule. We can then insert an n -dimensional vector in the header of every packet p , where the i^{th} element contains $o(p, \alpha_i)$ with α_i being the i^{th} hop in $\text{path}(p)$. Every time a packet arrives at a router, the router can pop the value at the head of this vector and use that as its priority (earlier values of output times get higher priority). This can perfectly replay any viable schedule (proof in Appendix A.2), which is not surprising, as having such detailed knowledge of the internal scheduling of the network is tantamount to knowing all the scheduling decisions made by the original algorithm. For reasons discussed previously, our definition limited the information available to the output time from the network as a whole, and not from each individual router; we call this *black-box* initialization.

Nonexistence of a UPS under black-box initialization: We can prove by counter-example (described in Appendix A.3) that *there is no UPS* under the conditions stated in §2.1. We provide some intuition for the counter-example later in this section. Given this impossibility result, we now ask *how close can we get to a UPS?*

Natural candidates for a near-UPS: Simple priority scheduling⁶ can reproduce all viable schedules on a single router, so it would seem to be a natural candidate for a near-UPS. However, for multihop networks it may be important to make the scheduling of a packet dependent on what has happened to it earlier in its path. For this, we consider Least Slack Time First (LSTF) [28].

In LSTF, each packet p carries its slack value in the packet header, which is initialized to $\text{slack}(p) = (o(p) - i(p) - t_{\min}(p, \text{src}(p), \text{dest}(p)))$ at the ingress; where $\text{src}(p)$ is the ingress of p , $\text{dest}(p)$ is the egress of p and $t_{\min}(p, \alpha, \beta)$ is the time p takes to go from router α to router β in an uncongested network. Therefore, the slack of a packet indicates the maximum queuing time (excluding the transmission time at any router) that the packet could tolerate without violating the replay condition. Each router, then, schedules the packet which has the

⁶By simple priority scheduling, we mean that the ingress assigns priority values to the packets and the routers simply schedule packets based on these static priority values.

least remaining slack at the time when its last bit is transmitted. Before forwarding the packet, the router overwrites the slack value in the packet’s header with its remaining slack (*i.e.*, the previous slack time minus the duration for which it waited in the queue before being transmitted).

An alternate way to implement this algorithm is having a static packet header as in Earliest Deadline First (EDF) and using additional state in the routers (reflecting the value of t_{\min}) to compute the priority for a packet at each router, but here we chose to use an approach with dynamic packet state. We provide more details about EDF and prove its equivalence to LSTF in Appendix A.5.

Key Results: Our analysis shows that the difficulty of replay is determined by the number of *congestion points*, where a *congestion point* is defined as a node where a packet is forced to “wait” during a given schedule.⁷ Our theorems show the following key results:

1. Priority scheduling can replay all viable schedules with no more than one congestion point per packet, and there are viable schedules with no more than two congestion points per packet that it cannot replay. (Proof in Appendix A.6.)
2. LSTF can replay all viable schedules with no more than two congestion points per packet, and there are viable schedules with no more than three congestion points per packet that it cannot replay. (Proof in Appendix A.7.)
3. There is no scheduling algorithm (obeying the aforementioned constraints on UPSs) that can replay *all* viable schedules with no more than three congestion points per packet, and the same holds for larger numbers of congestion points. (Proof in Appendix A.3.)

Main Takeaway: *LSTF is closer to being a UPS than simple priority scheduling, and no other candidate UPS can do better in terms of handling more congestion points.*

Intuition: It is clear why LSTF is superior to priority scheduling: by carrying information about previous delays in the packet header (in the form of the *remaining* slack value), LSTF can “make up for lost time” at later congestion points, whereas for priority scheduling packets with low priority might repeatedly get delayed (and thus miss their target output times).

We now provide some intuition for why LSTF works for two congestion points and not for three, by presenting an outline of the proof detailed in Appendix A.7. We define the *local deadline* of a packet p at a router α as the time when p is scheduled by α in the original schedule. The *global deadline* of p at α is defined as the time by when p must leave α

⁷For our theoretical results, we adopt a pessimistic definition of a congestion point, where a router that falls in the path of more than one flow is a congestion point (along with routers having output link capacity less than input link capacity or non work-conserving original schedules that make a packet wait explicitly). Since this definition is independent of per-packet dynamics, the set of congestion points remains the same in the original schedule and in the replay. This pessimistic definition is not required in practice, where the difficulty of replay would depend on the number of routers in a packet’s path which see significant queuing.

in order to meet its target output time, assuming that it sees no queuing delay after α . Hence, *global deadline* is the time when p 's slack at α becomes zero. We can prove that *as long as all packets arrive at a router at or before their local deadlines during the LSTF replay, no packet can miss its global deadline at α (i.e. no packet can have a negative slack at α)*. The proof for this follows from the fact that if all packets arrive at or before their *local deadline* at α , there exists a *feasible* schedule where no packet misses its *global deadline* at α (this *feasible* schedule is the same as the original schedule at α). We can now apply the standard LSTF (or EDF) optimality proof technique for a single processor [30], to show that this feasible schedule can be iteratively *transformed* to a feasible LSTF schedule at router α .

When there are only two congestion points per packet, it is guaranteed that every packet arrives at or before its local deadline at each congestion point during the LSTF replay. A packet can never arrive after its *local deadline* at its first congestion point, because it sees no queuing before that. Moreover, the *local deadline* is the same as the *global deadline* at the last congestion point. Therefore, if a packet arrives after its *local deadline* at its second (and last) congestion point, it means that it must have already missed its *global deadline* earlier, which, again, is not possible.

However, when there are three congestion points per packet, there is no guarantee that every packet arrives at or before its local deadline at each congestion point during the LSTF replay (due to the presence of a "middle" congestion point). One can, therefore, create counterexamples where unless LSTF (or, in fact, any other scheduling algorithm) makes precisely the right choice at the first congestion point of a packet p , at least one packet will miss its target output time, due to p arriving after its *local deadline* at its middle congestion point. We present such a counterexample in Appendix A.3, where we illustrate two ways of scheduling the same set of packets (having the same input times and paths) on a given topology with three congestion points per packet, resulting in two cases. The output times for two of the packets (named a and x), which compete with each other at the first congestion point (α_0), remains the same in both cases. However, one case requires scheduling a before x at α_0 and the second case requires scheduling x before a at α_0 , else a packet will end up missing its target output time at the second (or middle) congestion points of a and x respectively. Since the information available for header initialization for the two packets is the same in both cases, no deterministic scheduling algorithm with blackbox header initialization can make the *correct* choice at the first congestion point *in both cases*.

2.3 Empirical Results

The previous section clarified the theoretical limits on a *perfect* replay. Here we investigate, via ns-2 simulations [6], how well (a non-preemptable version of) LSTF

can *approximately* replay schedules in realistic networks.

Experiment Setup: Default scenario. We use a simplified Internet-2 topology [3], identical to the one used in [31] (consisting of 10 core routers connected by 16 links). We connect each core router to 10 edge routers using 1Gbps links and each edge router is attached to an end host via a 10Gbps link. The number of hops per packet is in the range of 4 to 7, excluding the end hosts. We refer to this topology as I2 1Gbps-10Gbps. Each end host generates UDP flows using a Poisson inter-arrival model, with the destination picked randomly for each flow. Our default scenario runs at 70% utilization. The flow sizes are picked from a heavy-tailed distribution [11, 12]. Since our focus is on packet scheduling, not dropping policies, we use large buffer sizes that ensure no packet drops. Note that we use higher than usual access bandwidths for our default scenario to increase the stress on the schedulers in the core routers, where the number of congestion points seen by most packets is two, three or four for 22%, 44% and 24% packets respectively.⁸ We also present results for smaller (and more realistic) access bandwidths, where most packets see smaller number of congestion points (one, two or three for 18%, 46% and 26% packets respectively), resulting in better replay performance.

Varying parameters. We tested a wide range of experimental scenarios by varying different parameters from their default values. We present results for a small subset of these scenarios here: (1) the default scenario with network utilization varied from 10-90% (2) the default scenario but with 1Gbps link between the endhosts and the edge routers (I2 1Gbps-1Gbps), with 10Gbps links between the edge routers and the core (I2 10Gbps-10Gbps) and with all link capacities in the I2 1Gbps-1Gbps topology reduced by a factor of 10 (I2 / 10) and (3) the default scenario applied to two different topologies, a bigger Rocketfuel topology [43] (with 83 core routers connected by 131 links) and a full bisection bandwidth datacenter (fat-tree) topology from [10] (with 10Gbps links). Note that our other results were generally consistent with those presented here.

Scheduling algorithms. Our default case, which we expected to be hard to replay, uses completely arbitrary schedules produced by a *random* scheduler (which picks the packet to be scheduled randomly from the set of queued up packets). We also present results for more traditional packet scheduling algorithms: FIFO, LIFO, fair queuing [20], and SJF (shortest job first using priorities). We also looked at two scenarios with a mixture of scheduling algorithms: one where half of the routers run FIFO+ [19] and the other half run fair queuing, and one where fair queuing is used to isolate two classes of traffic, with one class being scheduled with SJF and the

⁸To compute this, we record the number of non-empty queues (excluding the endhost queues) encountered by each packet.

Topology	Avg. Link Utilization	Scheduling Algorithm	Fraction of packets overdue Total > T	
I2 1Gbps-10Gbps	70%	Random	0.0021	0.0002
I2 1Gbps-10Gbps	10% 30% 50% 90%	Random	0.0007 0.0281 0.0221 0.0008	0.0 0.0017 0.0002 4×10^{-6}
I2 1Gbps-1Gbps I2 10Gbps-10Gbps I2 / 10	70%	Random	0.0204 0.0631 0.0127	8×10^{-6} 0.0448 0.00001
Rocketfuel Datacenter	70%	Random	0.0246 0.0164	0.0063 0.0154
I2 1Gbps-10Gbps	70%	FIFO FQ SJF LIFO FQ/FIFO+ FQ: SJF/FIFO	0.0143 0.0271 0.1833 0.1477 0.0152 0.0297	0.0006 0.0002 0.0019 0.0067 0.0004 0.0003

Table 1: LSTF replay performance across various scenarios. T represents the transmission time at the bottleneck link.

other class being scheduled with FIFO.

Evaluation Metrics: We consider two metrics. First, we measure the fraction of packets that are overdue (i.e., which do not meet the original schedule’s target). Second, to capture the *extent* to which packets fail to meet their targets, we measure the fraction of packets that are overdue by more than a threshold value T , where T is one transmission time on the bottleneck link ($\approx 12\mu s$ for 1Gbps). We pick this value of T both because it is sufficiently small that we can assume being overdue by this small amount is of negligible practical importance, and also because this is the order of violation we should expect given that our implementation of LSTF is non-preemptive. While we may have many small violations of replay (because of non-preemption), one would hope that most such violations are less than T .

Results: Table 1 shows the simulation results for LSTF replay for various scenarios, which we now discuss.

(1) Replayability. Consider the column showing the fraction of packets overdue. In all but three cases (we examine these shortly) over 97% of packets meet their target output times. In addition, the fraction of packets that did not arrive within T of their target output times is much smaller; even in the worst case of SJF scheduling (where 18.33% of packets failed to arrive by their target output times), only 0.19% of packets are overdue by more than T . Most scenarios perform substantially better: e.g., in our default scenario with Random scheduling, only 0.21% of packets miss their targets and only 0.02% are overdue by more than T . Hence, we conclude that even without preemption LSTF achieves good (but not perfect) replayability under a wide range of scenarios.

(2) Effect of varying network utilization. The second row in Table 1 shows the effect of varying network utilization. We see that at low utilization (10%), LSTF achieves exceptionally good replayability with a total of

only 0.07% of packets overdue. Replayability deteriorates as utilization is increased to 30% but then (somewhat surprisingly) improves again as utilization increases. This improvement occurs because with increasing utilization, the amount of queuing (and thus the average slack across packets) in the original schedule also increases. This provides more room for slack re-adjustments when packets wait longer at queues seen early in their paths during the replay. We observed this trend in all our experiments though the exact location of the “low point” varied across settings.

(3) Effect of varying link bandwidths. The third row shows the effect of changing the relative values of access/edge vs. core links. We see that while decreasing access link bandwidth (I2 1Gbps-1Gbps) resulted in a much smaller fraction of packets being overdue by more than T (0.0008%), increasing the edge-to-core link bandwidth (I2 10Gbps-10Gbps) resulted in a significantly higher fraction (4.48%). For I2 1Gbps-1Gbps, packets are paced by the endhost link, resulting in few congestion points thus improving LSTF’s replayability. In contrast, with I2 10Gbps-10Gbps, both the access and edge links have a higher bandwidth than most core links; hence packets (that are no longer paced at the endhosts or the edges) arrive at the core routers very close to one another and hence the effect of one packet being overdue *cascades* over to the following packets. Decreasing the absolute bandwidths in I2 / 10, while keeping the ratio between access and edge links the same as that in I2 1Gbps-1Gbps, did not produce significantly different results compared to I2 1Gbps-1Gbps, indicating that the relative link capacities have a greater impact on the replay performance than the absolute link capacities.

(4) Effect of varying topology. The fourth row in Table 1 shows our results using different topologies. LSTF performs well in both cases: only 2.46% (Rocketfuel) and 1.64% (datacenter) of packets fail replay. These numbers are still somewhat higher than our default case. The reason for this is similar to that for the I2 10Gbps-10Gbps topology – all links in the datacenter fat-tree topology are set to 10Gbps, while in our simulations, we set half of the core links in the Rocketfuel topology to have bandwidths smaller than the access links.

(5) Varying Scheduling Algorithms. Row five in Table 1 shows LSTF’s ability to replay different scheduling algorithms. We see that LSTF performs well for FIFO, FQ, and the combination cases (a mixture of FQ/FIFO+ and having FQ share between FIFO and SJF); e.g., with FIFO, fewer than 0.06% of packets are overdue by more than T . However, there are two problematic cases: SJF and LIFO fare worse with 18.33% and 14.77% of packets failing replay (although only 0.19% and 0.67% of packets are overdue by more than T respectively). The reason stems from a combination of two factors: (1) for these algorithms a larger fraction of packets have a very small slack value (as one might expect from the scheduling logic which

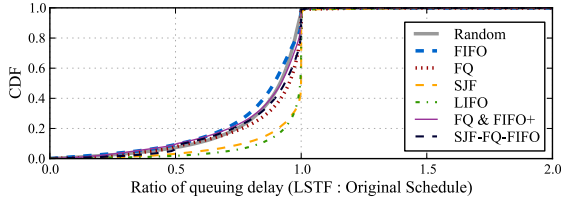


Figure 1: Ratio of queuing delay with varying packet scheduling algorithms, on I2 1Gbps-10Gbps topology at 70% utilization.

produces a larger skew in the slack distribution), and (2) for these packets with small slack values, LSTF *without preemption* is often unable to “compensate” for misspent slack that occurred earlier in the path. To verify this intuition, we extended our simulator to support preemption and repeated our experiments: with preemption, the fraction of packets that failed replay dropped to 0.24% (from 18.33%) for SJF and to 0.25% (from 14.77%) for LIFO.

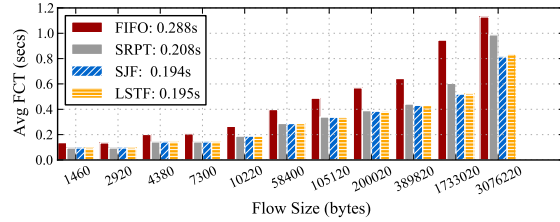
(6) End-to-end (Queuing) Delay. Our results so far evaluate LSTF in terms of measures that we introduced to test universality. We now evaluate LSTF using the more traditional metric of packet delay, focusing on the queuing delay a packet experiences. Figure 1 shows the CDF of the ratios of the queuing delay that a packet sees with LSTF to the queuing delay that it sees in the original schedule, for varying packet scheduling algorithms. We were surprised to see that most of the packets actually have a smaller queuing delay in the LSTF replay than in the original schedule. This is because LSTF eliminates “wasted waiting”, in that it never makes packet A wait behind packet B if packet B is going to have significantly more waiting later in its path.

(7) Comparison with Priorities. To provide a point of comparison, we also did a replay using simple priorities for our default scenario, where the priority for a packet p is set to $o(p)$ (which seemed most intuitive to us). As expected, the resulting replay performance is much worse than LSTF: 21% packets are overdue in total, with 20.69% being overdue by more than T . For the same scenario, LSTF has only 0.21% packets overdue in total, with merely 0.02% packets overdue by more than T .

Summary: We observe that, in almost all cases, less than 1% of the packets are overdue with LSTF by more than T . The replay performance initially degrades and then starts improving as the network utilization increases. The distribution of link speeds has a bigger influence on the replay results than the scale of the topology. Replay performance is better for scheduling algorithms that produce a smaller skew in the slack distribution. LSTF replay performance is significantly better than simple priorities replay performance, with the most intuitive priority assignment.

3 Practical: Achieving Various Objectives

While replayability demonstrates the theoretical flexibility of LSTF, it does not provide evidence that it would be practically useful. In this section we look at how LSTF



Expt. Setup	Avg FCT (s)			
	FIFO	SRPT	SJF	LSTF
I2 1Gbps-10Gbps at 30% util.	0.189	0.183	0.182	0.182
I2 1Gbps-10Gbps at 50% util.	0.212	0.189	0.185	0.185
I2 1Gbps-10Gbps at 70% util.	0.288	0.208	0.194	0.195
I2 1Gbps-1Gbps at 70% util.	0.252	0.209	0.202	0.202
I2 / 10 at 70% util.	0.899	0.658	0.620	0.621
Rocketfuel at 70% util.	0.305	0.240	0.228	0.228
Datacenter at 70% util.	0.058	0.018	0.016	0.015

Figure 2: The graph shows the average FCT bucketed by flow size obtained with FIFO, SRPT and SJF (using priorities and LSTF) for I2 1Gbps-10Gbps at 70% utilization. The legend indicates the average FCT across all flows. The table indicates the average FCTs for varying settings.

can be used *in practice* to meet the following performance objectives: minimizing average flow completion times, minimizing tail latencies, and achieving per-flow fairness.

Since the knowledge of a previous schedule is unavailable in practice, instead of using a given set of output times (as done in §2.3), we now use heuristics to assign the slacks in an effort to achieve these objectives. Our goal here is not to outperform the state-of-the-art for each objective in all scenarios, but instead we aim to be competitive with the state-of-the-art in most common cases.

In presenting our results for each objective, we first describe the slack initialization heuristic we use and then present some ns-2 [6] simulation results on (i) how LSTF performs relative to the state-of-the-art scheduling algorithm and (ii) how they both compare to FIFO scheduling (as a baseline to indicate the overall impact of specialized scheduling for this objective). As our default case, we use the I2 1Gbps-10Gbps topology using the same workload as in the previous section (running at 70% average utilization). We also present aggregate results at different utilization levels and for variations in the default topology (I2 1Gbps-1Gbps and I2 / 10), for the bigger Rocketfuel topology, and for the datacenter topology (for selected objectives). The switches use non-preemptive scheduling (including for LSTF) and have finite buffers (packets with the highest slack are dropped when the buffer is full). Unless otherwise specified, our experiments use TCP flows with router buffer sizes of 5MB for the WAN simulations (equal to the average bandwidth-delay product for our default topology) and 500KB for the datacenter simulations.

3.1 Average Flow Completion Time

While there have been several proposals on how to minimize flow completion time (FCT) via the transport protocol [21, 31], here we focus on *scheduling*’s impact on FCT, while using standard TCP New Reno at the endpoints. In [10] it is shown that (i) Shortest Remaining Processing

Time (SRPT) is close to optimal for minimizing the mean FCT and (ii) Shortest Job First (SJF) produces results similar to SRPT for realistic heavy-tailed distribution. Thus, these are the two algorithms we use as benchmarks.

Slack Initialization: We make LSTF emulate SJF by initializing the slack for a packet p as $slack(p) = fs(p) * D$, where $fs(p)$ is the size of the flow to which p belongs (in terms of the number of MSS-sized packets in the flow) and D is a value much larger than the queuing delay seen by any packet in the network. We use a value of $D = 1$ sec for our simulations.

Evaluation: Figure 2 compares LSTF with three other scheduling algorithms – FIFO, SJF and SRPT with *starvation prevention* as in [10]. Both SJF and SRPT have significantly lower mean FCT than FIFO. The LSTF based execution of SJF produces nearly the same results as the strict priorities based execution.

We also look at how in-network scheduling can be used along with changes in the endhost TCP stack to achieve the same objective in Appendix B.

3.2 Tail Packet Delays

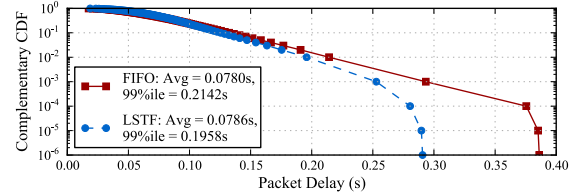
Clark et. al. [19] proposed the FIFO+ algorithm, where packets are prioritized at a router based on the amount of queuing delay they have seen at their previous hops, for minimizing the tail packet delays in multi-hop networks. FIFO+ is *identical* to LSTF scheduling where all packets are initialized with the same slack value.

Slack Initialization: All incoming packets are initialized with the same slack value (we use an initial slack value of 1 second in our simulations). With the slack update taking place at every router, the packets that have waited longer in the network queues are naturally given preference over those that have waited for a smaller duration.

Evaluation: We compare LSTF (which, with the above slack initialization, is identical to FIFO+) with FIFO, the primary metric being the 99%ile end-to-end one way delay seen by the packets. Figure 3 shows our results. To better understand the impact of the two scheduling policies on the packet delays, our evaluation uses an open-loop setting with UDP flows. With LSTF, packets that have traversed through more number of hops, and have therefore spent more slack in the network, get preference over shorter-RTT packets that have traversed through fewer hops. While this might produce a slight increase in the average packet delay, it reduces the tail. This is in-line with the observations made in [19].

3.3 Fairness

Fairness is a common scheduling goal, which involves two different aspects: *asymptotic* bandwidth allocation (eventual convergence to the fair-share rate) and *instantaneous* bandwidth allocation (enforcing this fairness on small



Expt. Setup	Avg Delay (s)		99%ile Delay (s)	
	FIFO	LSTF	FIFO	LSTF
I2 1Gbps-10Gbps at 30% util.	0.0411	0.0411	0.0911	0.0868
I2 1Gbps-10Gbps at 50% util.	0.0516	0.0517	0.1288	0.1195
I2 1Gbps-10Gbps at 70% util.	0.0780	0.0786	0.2142	0.1958
I2 1Gbps-1Gbps at 70% util.	0.0771	0.0771	0.2163	0.216
I2 / 10 at 70% util.	0.5762	0.5765	1.9393	1.9367
Rocketfuel at 70% util.	0.1891	0.1883	3.8139	3.7199
Datacenter at 70% util.	0.0250	0.0240	0.1352	0.1100

Figure 3: Tail packet delays for LSTF compared to FIFO. The graph shows the complementary CDF of packet delays for the I2 1Gbps-10Gbps topology at 70% utilization with the average and 99%ile packet delay values indicated in the legend. The table shows the corresponding results for varying settings.

time-scales, so every flow experiences the equivalent of a per-flow pipe). The former can be measured by looking at long-term throughput measures, while the latter is best measured in terms of the flow completion times of relatively short flows (which measures bandwidth allocation on short time scales). We now show how LSTF can be used to achieve both of these goals, but more effectively the former than the latter. Our slack assignment heuristic can also be easily extended to achieve weighted fair queuing, but we do not present those results here.

Slack Initialization: The slack assignment for fairness works on the assumption that we have some ballpark notion of the fair-share rate for each flow and that it does not fluctuate wildly with time. Our approach to assigning slacks is inspired from [46]. We assign $slack = 0$ to the first packet of the flow and the slack of any subsequent packet p_i is then initialized as:

$$slack(p_i) = \max\left(0, slack(p_{i-1}) + \frac{size(p_i)}{r_{est}} - (i(p_i) - i(p_{i-1}))\right)$$

where $i(p)$ is the arrival time of a packet p at the ingress, $size(p)$ is its size in bits, and r_{est} is an estimate of the fair-share rate r^* in bps. We show that the above heuristic leads to asymptotic fairness, for *any* value of r_{est} that is less than r^* , as long as all flows use the same value. The same heuristic can also be used to provide instantaneous fairness, when we have a complex mix of short-lived flows, where the r_{est} value that performs the best depends on the link bandwidths and their utilization levels. A reasonable value of r_{est} can be estimated using knowledge about the network topology and traffic matrices, though we leave a detailed exploration of this to future work.

Evaluation: Asymptotic Fairness. We evaluate the asymptotic fairness property by running our simulation on the Internet2 topology with 10Gbps edges, such that all the congestion happens at the core. However, we reduce the propagation delay to $10\mu s$ for each link, to make

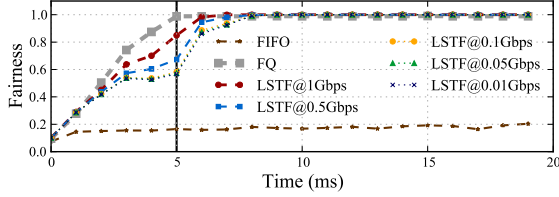


Figure 4: Fairness for long-lived flows on Internet2 topology. The legend indicates the value of r_{est} used for LSTF slack initialization.

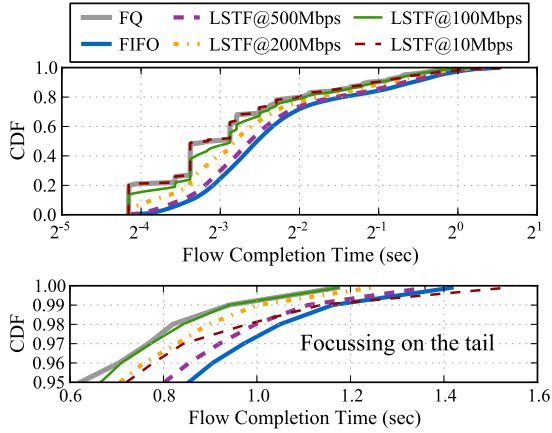


Figure 5: CDF of FCTs for the I2 1Gbps-10Gbps topology at 70% utilization.

Expt. Setup	Avg FCT across bytes (s)			Best r_{est} (Mbps)	Reasonable r_{est} Range (Mbps)
	FIFO	FQ	LSTF		
I2 1Gbps-10Gbps at 30% util.	0.563	0.537	0.538	300	10-900
I2 1Gbps-10Gbps at 50% util.	0.626	0.549	0.555	200	10-800
I2 1Gbps-10Gbps at 70% util.	0.811	0.622	0.632	100	50-200
I2 1Gbps-1Gbps at 70% util.	0.766	0.630	0.652	100	50-400
I2 / 10 at 70% util.	4.838	2.295	2.759	10	10-20
Rocketfuel at 70% util.	0.964	0.796	0.824	100	50-300

Table 2: FCT averaged across bytes for FIFO, FQ and LSTF (with best r_{est} value) across varying settings. The last column indicates the range of r_{est} values that produce results within 10% of the best r_{est} result.

the experiment more scalable, while the buffer size is kept large (50MB) so that fairness is dominated by the scheduling policy and not by how TCP reacts to packet drops. We start 90 long-lived TCP flows with a random jitter in the start times ranging from 0-5ms. The topology is such that the fair share rate of each flow on each link in the core network (which is shared by up to 13 flows) is around 1Gbps. We use different values for $r_{est} \leq 1$ Gbps for computing the initial slacks and compare our results with fair queuing (FQ). Figure 4 shows the fairness computed using Jain’s Fairness Index [27], from the throughput each flow receives per millisecond. Since we use the throughput received by each of the 90 flows to compute the fairness index, it reaches 1 with FQ only at 5ms, after all the flows have started. We see that LSTF is able to converge to perfect fairness, even when r_{est} is 100X smaller than r^* . It converges slightly sooner when r_{est} is closer to r^* , though the subsequent differences in the time to convergence decrease with decreasing values of r_{est} .

The detailed explanation of how this works along with

more evaluation (on multiple bottlenecks and weighted fairness) has been provided in Appendix C.

Evaluation: Instantaneous Fairness. As one might expect, the choice of r_{est} has a bigger impact on instantaneous fairness than on asymptotic fairness. A very high r_{est} value would not provide sufficient isolation across flows. On the other hand, a very small r_{est} value can starve the long flows. This is because the assigned slack values for the later packets of long flows with high sequence numbers would be much higher than the actual slack they experience. As a result, they will end up waiting longer in the queues, while the initial packets of newer flows with smaller slack values would end up getting a higher precedence.

To verify this intuition, we evaluated our LSTF slack assignment scheme by running our standard workload with a mix of TCP flows ranging from sizes 1.5KB - 3MB on our default I2 1Gbps-10Gbps topology at 70% utilization, with 50MB buffer size. Note that the traffic pattern is now bursty and the instantaneous utilization of a link is often lower or higher than the assigned average utilization level. The CDF of the FCTs thus obtained is shown in Figure 5. As expected, the distribution of FCTs looks very different between FQ and FIFO. FQ isolates the flows from each-other, significantly reducing the FCT seen by short to medium size flows, compared to FIFO. The long flows are also helped a little by FQ, again due to the isolation provided from one-another.

LSTF performance varies somewhere in between FIFO and FQ, as we vary r_{est} values between 500Mbps to 10Mbps. A high value of $r_{est} = 500$ Mbps does not provide sufficient isolation and the performance is close to FIFO. As we reduce the value of r_{est} , the “isolation-effect” increases. However, for very small r_{est} values (e.g. 10Mbps), the tail FCT (for the long flows) is much higher than FQ, due to the starvation effect explained before.

We try to capture this trade-off between isolation for short and medium sized flows and starvation for long flows, by using average FCT across bytes (in other words, the average FCT weighted by flow size) as our key metric. We term the r_{est} value that achieves the sweetest spot in this trade-off as the “best” r_{est} value. The r_{est} values that produce average FCT which is within 10% of the value produced by the best r_{est} are termed as “reasonable” r_{est} values. Table 2 presents our results across different settings. We find that (1) LSTF produces significantly lower average FCT than FIFO, performing only slightly worse than FQ (2) As expected, the best r_{est} value decreases with increasing utilization and with decreasing bandwidths (as in the case of I2 / 10 topology), while the range of reasonable r_{est} values gets narrower with increasing utilization and with decreasing bandwidths.

Thus, for instantaneous fairness, LSTF would require some estimate of the per-flow rate. We believe that this can be obtained from the knowledge of the network topology

(in particular, the link bandwidths), which is available to the ISPs, and on-line measurement of traffic matrices and link utilization levels, which can be done using various tools [14, 18, 35]. However, this does impose a higher burden on deploying LSTF than on FQ or other such scheduling algorithms.

3.4 Limitations of LSTF: Policy-based objectives

So far we showed how LSTF achieves various performance objectives. We now describe certain policy-based objectives that are hard to achieve with LSTF.

Multi-tenancy: As network virtualization becomes more popular, networks are often called upon to support multiple tenants or traffic classes, with each having their own networking objectives. Network providers can enforce isolation across such tenants (or classes of traffic) through static bandwidth provisioning, which can be implemented via dedicated hard-wired links [1, 5] or through multiqueue scheduling algorithms such as fair queuing or round robin [20]. LSTF can work in conjunction with both of these isolation mechanisms to meet different desired performance objectives for each tenant (or class of traffic).

However, without such multiqueue support it cannot provide such isolation or fairness on a per-class or per-tenant basis. This is because for class-based fairness (which also includes hierarchical fairness) the appropriate slack assignment for a packet at a particular ingress depends on the input from other ingresses (since these packets can belong to the same class). Note, however, that if two or more classes/tenants are separated by strict prioritization, LSTF can be used to enforce the appropriate precedence order, along with meeting the individual performance objective for each class.

Traffic Shaping: Shaping or rate limiting flows at a particular router requires non-work-conserving algorithms such as Token Bucket Filters [8]. LSTF itself is a work-conserving algorithm and cannot shape or rate limit the traffic on its own. We believe that shaping the traffic only at the edge, with the core remaining work-conserving, would also produce the desired network-wide behavior, though this requires further exploration.

4 Incorporating Network Feedback

Up until now we have considered packet scheduling in isolation, whereas in the Internet today routers send implicit feedback to hosts via packet drops [22, 32] (or marking, as in ECN [37]). This is often called Active Queue Management (AQM), and its goal is to reduce the per-packet delays while keeping throughput high. We now consider how we might generalize our LSTF approach to incorporate such network feedback as embodied in AQM schemes.

LSTF is just a scheduling algorithm and cannot perform AQM on its own. Thus, at first glance, one might think that incorporating AQM into LSTF would require

implementing the AQM scheme in each router, which would then require us to find a universal AQM scheme in order to fulfill our pursuit of universality. On the contrary, LSTF enables a novel edge-based approach to AQM based on the following insights: (1) As long as appropriate packets are chosen, it does not matter *where* they are being dropped (or marked) – whether it is inside the core routers or at the edge. (2) In addition to scheduling packets LSTF produces a very useful by-product, carried by the slack values in the packets, which gives us a precise measure of the one-way queuing delay seen by the packet and can be used for AQM. For obtaining this by-product, an extra field is added to the packet header at the ingress which stores the assigned slack value (called the initial slack field), which remains untouched as the packet traverses the network. The other field where the ingress stores the assigned slack value is updated as per the LSTF algorithm; we call this the current slack field. The precise amount of queuing delay seen by the packet within the network (or the *used slack* value) can be computed at the edge by simply comparing the initial slack field and the current slack field.

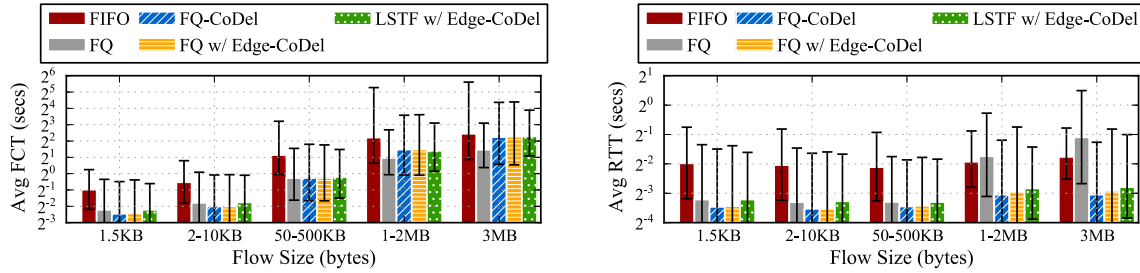
We evaluate our edge-based approach to AQM in the context of (1) CoDel [32], the state-of-the-art AQM scheme for wide area networks and (2) ECN used with DCTCP [9], the state-of-the-art AQM scheme for datacenters.

4.1 Emulating CoDel from Edge

Background: In CoDel, the amount of time a packet has spent in a queue is recorded as the sojourn time. A packet is dropped if its sojourn time exceeds a fixed target (set to 5ms [33]), and if the last packet drop happened beyond a certain interval (initialized to 100ms [33]). When a packet is dropped, the interval value is reduced using a control law, which divides the initial interval value by the square root of the number of packets dropped. The interval is refreshed (re-initialized to 100ms) when the queue becomes empty, or when a packet sees a sojourn time less than the target.⁹ An extension to CoDel is FQ-CoDel [25], where the scheduler round-robins across flows and the CoDel control loop is applied to each flow individually. The interval for a flow is refreshed when there are no more packets belonging to that flow in the queue. FQ-CoDel is considered to be better than CoDel in all regards, even by one of the co-developers of CoDel [4].

Edge-CoDel: We aim to approximate FQ-CoDel from the edge by using LSTF to implement per-flow fairness in routers (as in §3.3). We then compute the used slack value at the egress router for every packet, as described above, and run the FQ-CoDel logic for when to drop packets for each flow, keeping the control law and the parameters (the target value and the initial interval value) the same as in

⁹CoDel is a little more complicated than this, and while our implementation follows the CoDel specification [33], our explanation has been simplified, highlighting only the relevant points for brevity.



Expt. Setup	r_{est} (Mbps)	Avg FCT across bytes (s)					Avg RTT across bytes (s)				
		FIFO	FQ	FQ-CoDel	FQ w/ Edge-CoDel	LSTF w/ Edge-CoDel	FIFO	FQ	FQ-CoDel	FQ w/ Edge-CoDel	LSTF w/ Edge-CoDel
I2 1Gbps-10Gbps at 70% util.	100	0.811	0.622	0.642	0.633	0.641	0.0756	0.0733	0.0642	0.0646	0.0661
I2 1Gbps-1Gbps at 70% util.	100	0.766	0.630	0.642	0.637	0.658	0.0716	0.0702	0.0639	0.0643	0.0666
I2 / 10 at 30% util.	40	0.918	0.836	0.897	0.887	0.907	0.0998	0.1085	0.0792	0.0798	0.0826
I2 / 10 at 50% util.	30	1.706	1.214	1.430	1.369	1.427	0.1384	0.1752	0.0901	0.0918	0.1001
I2 / 10 at 70% util.	10	4.837	2.295	3.687	3.738	3.739	0.2779	0.3752	0.1182	0.1281	0.1388
I2 / 10, half RTTs at 70% util.	10	4.569	2.023	3.196	3.245	3.405	0.2555	0.3607	0.0995	0.1131	0.1165
I2 / 10, double RTTs at 70% util.	10	5.098	2.769	4.243	4.125	4.389	0.325	0.4172	0.1591	0.1640	0.1843
Rocketfuel at 70% util.	100	0.964	0.796	0.840	0.813	0.835	0.0922	0.0991	0.0794	0.0788	0.0836

Figure 6: The figures show the average FCT and RTT values for I2 / 10 at 70% utilization (LSTF uses fairness slack assignment with $r_{est} = 10Mbps$). The error bars indicate the 10th and the 99th percentile values and the y-axis is in log-scale. The table indicates the average FCT and RTTs (across bytes) for varying settings.

FQ-CoDel. We call this approach Edge-CoDel.

There are only two things that change in Edge-CoDel as compared to FQ-CoDel. First, instead of looking at the sojourn time of each queue individually, Edge-CoDel looks at the total queuing time of the packet across the entire network. The second change is with respect to how the CoDel interval is refreshed. As mentioned before, in traditional FQ-CoDel, there are two events that trigger a refresh in the interval (i) when a packet’s sojourn time is less than the target and (ii) when all the queued-up packets for a given flow have been transmitted. While Edge-CoDel can react to the former, it has no explicit way of knowing the latter. To address this, we refresh the interval if the difference in the send time of two consecutive packets (found using TCP timestamps that are enabled by default) is more than a certain threshold. Clearly, this *refresh threshold* must be greater than CoDel’s target queuing delay value. We find that a refresh threshold of 2-4 times the target value (10-20ms) works reasonably well.

Evaluation: In our experiments, we compare four different schemes: (1) FIFO without AQM (to set a baseline), (2) FQ without AQM (to see the effects of FQ on its own), (3) FQ-CoDel (to provide the state-of-the-art comparison) (4) LSTF scheduling (with slacks assigned to meet the fairness objective using appropriate r_{est} values) in conjunction with Edge-CoDel. As we move from (3) to (4), we make two transitions – first is with respect to the scheduling done inside the network (perfect isolation with FQ vs approximate isolation with LSTF) and the second is the shift of AQM logic from inside the network to the edge. Therefore, as an incremental step in between the two transitions, we also provide results for FQ with Edge-CoDel, where routers do FQ across flows (with the slack values maintained only

for book-keeping) and AQM is done by Edge-CoDel. This allows us see how well Edge-CoDel works with perfect per-router isolation. The refresh threshold we use for Edge-CoDel in both cases is 20ms (4 times the CoDel target value). The buffer size is increased to 50MB so that AQM kicks in before a natural packet drop occurs.

Figure 6 shows our results for varying settings and schemes. The main metrics we use for evaluation are the FCTs and the per-packet RTTs, since the goal of an AQM scheme is to maintain high throughput (or small FCTs) while keeping the RTTs small. The two graphs show the average FCT and the average RTT across flows bucketed by their size for the I2 / 10 topology at 70% utilization (where AQM produces a bigger impact compared to our default case). As expected, we find that while FQ helps in reducing the FCT values as compared to FIFO, it results in significantly higher RTTs than FIFO for long flows. FQ-CoDel reduces the RTT seen by long flows compared to FQ (with the short flows having RTT smaller than FIFO and comparable to FQ). What is new is that, shifting the CoDel logic to the edge through Edge-CoDel while doing FQ in the router makes very little difference as compared to FQ-CoDel. As we experiment with varying settings, we find that in some cases, FQ with Edge-CoDel results in slightly smaller FCTs at the cost of slightly higher RTTs than FQ-CoDel. We believe that this is due to the difference in how the CoDel interval is refreshed with Edge-CoDel and with in-router FQ-CoDel. Replacing the scheduling algorithm with LSTF again produces minor differences in the results compared to FQ-CoDel. Both the FCT and the RTT are slightly higher than FQ-CoDel for almost all cases, and we attribute the differences to LSTF’s *approximation* of round-robin service across flows. Nonetheless, the average FCTs obtained are significantly

Util.	Avg FCT (s)			Avg RTT (ms)		
	FIFO w/ No ECN	FIFO w/ ECN	LSTF w/ Edge- ECN	FIFO w/ No ECN	FIFO w/ ECN	LSTF w/ Edge- ECN
30%	0.0020	0.0011	0.0011	0.2069	0.1123	0.1077
50%	0.0219	0.0086	0.0079	0.3425	0.1601	0.1477
70%	0.0501	0.0241	0.0240	0.4497	0.2616	0.2494

Table 3: DCTCP performance with no ECN, ECN (in-switch) and Edge-ECN for the datacenter topology at varying utilizations.

lower than FIFO and the average RTTs are significantly lower than both FIFO and FQ for all cases.

Varying the refresh threshold used for Edge-CoDel produces minor differences in the aggregate results, a detailed evaluation of which can be found in Appendix D.

4.2 Emulating ECN for DCTCP from Edge

Background: DCTCP [9] is a congestion control scheme for datacenters that uses ECN marks as a congestion signal to control the queue size before a packet drop occurs. It requires the routers to mark the packets whenever the instantaneous queue size goes beyond a certain threshold K . These markings are echoed back to the sender with the acknowledgments and the sender decreases its sending rate in proportion to the ECN marked packets.

Edge-ECN: The marking process can be moved to the edge (or the receiving endhost) by simply marking a packet if its queuing delay (computed, as in §4.1, by subtracting the initial slack value from the current slack value) is greater than the transmission time of K packets. This transmission time is easy to compute in datacenters where the link capacities are known.

Evaluation: The results for varying utilization levels are shown in Table 3. We compare Edge-ECN running LSTF in the routers (with all packets initialized to the same slack value) with in-switch ECN running FIFO in the routers, both using the same unmodified DCTCP algorithm at the endhosts. We use the DCTCP default value of $K = 15$ packets as the marking threshold. We also present results for DCTCP with no ECN marks (which reduces to TCP) and FIFO scheduling, as a comparison point. We see that both in-switch ECN and Edge-ECN DCTCP have comparable performance, with significantly lower average FCTs and RTTs than no ECN TCP.

Summary: The *used slack* information available as a by-product from LSTF can be effectively used to emulate an AQM scheme from the edge of the network.

5 LSTF Implementation

In this section, we study the feasibility of implementing LSTF in the routers. We start with showing that given a switch that supports fine-grained priority scheduling, it is trivial to implement LSTF on it using programmable header processing mechanisms [14, 15]. We then explore two different proposals for implementing fine-grained priorities in hardware.

Using fine-grained priorities to implement LSTF: Consider a packet p that arrives at a router α at time $i(p, \alpha)$, with slack $slack(p, \alpha)$. As mentioned in §2, LSTF prioritizes packets based on their remaining slack value at the time when their last bit is transmitted. This term is given by $(slack(p, \alpha) - (t - i(p, \alpha)) + T(p, \alpha))$ at any time t while p is waiting at α . $T(p, \alpha)$ is the transmission time of p at α , which is added to account for the remaining slack of p , relative to other packets, when its *last bit* is transmitted. Since t is same for all packets at any given point of time when the packets are being compared at α , the deciding term is $(slack(p, \alpha) + i(p, \alpha) + T(p, \alpha))$. With $slack(p, \alpha)$ being available in the packet header and the values of $i(p, \alpha)$ and $T(p, \alpha)$ being available at α when the packet arrives at the router, this term can be easily computed and attached to the packet as its priority value. Right before a packet p is transmitted by the router, its slack can be overwritten by the remaining slack value, computed by subtracting the stored priority value $(slack(p, \alpha) + i(p, \alpha) + T(p, \alpha))$ with the sum of the current time and $T(p, \alpha)$. We verified that these steps can be easily executed using P4 [14].

Implementing fine-grained priorities in hardware: Fine-grained priorities can be implemented by using specialized data-structures such as pipelined heap (p-heap) [13, 26], which can scale to very large buffers (>100MB), because the pipeline stage time is not affected by the queue size. However, p-heaps are difficult to implement and verify due to their intricate design and large chip area, thus resulting in higher costs. The p-heap implemented by Ioannou et. al. [26] using a 130nm technology node has a per-port area overhead of 10% (over a typical switching chip with minimum area of $200mm^2$ [23])¹⁰.

Leveraging the advancement in hardware technology over the years, Sivaraman et. al. [41] propose a simpler solution, based on bucket-sort algorithm. The area overhead reduces to only 1.65% (over a baseline single-chip shared-memory switch such as the Broadcom Trident [2]), when implemented using a 16nm technology node. While this approach is much cheaper to implement, it cannot scale to very large buffer sizes (beyond a few tens of MBs).

Thus, given these choices, it does not appear a significant challenge to implement LSTF at linespeed, though the key trade-offs between cost, simplicity and buffer limits need to be taken into consideration. To support a scale-out infrastructure, most datacenters today use a large number of inexpensive single chip shared memory switches [40], which have shallow buffers (around 10MB). The low overhead bucket-sort based approach [41] towards implementing LSTF would be ideal in such a setup. Core routers in wide area, on the other hand, have deep buffers (a few hundred MBs) and would require the more expensive

¹⁰130nm technology node was developed in 2001; the overheads would be lower for an implementation using the latest technology (14nm).

p-heap based implementation [13, 26]. While they are fewer in number [29], they may cost up to millions of dollars. Supporting the slightly more expensive, but flexible LSTF implementation would, to a large extent, obviate the need for replacing these expensive routers with changing demands, resulting in long-term savings. We are also optimistic that advancements in hardware technology would further reduce the cost overheads of implementing LSTF.

6 Related Work

The literature on packet scheduling is vast. Here we only touch on a few topics most relevant to our work.

The real-time scheduling literature has studied the optimality of scheduling algorithms¹¹ (in particular EDF and LSTF) for single and multiple processors [28, 30]. Liu and Layland [30] proved the optimality of EDF for a single processor in hard real-time systems. LSTF was then shown to be optimal for single-processor scheduling, while being more effective than EDF (though not optimal) for multi-processor scheduling [28]. In the context of networking, [17] provides theoretical results on emulating the schedules produced by a single output-queued switch using a combined input-output queued switch with a smaller speed-up of at most two. To the best of our knowledge, the optimality or universality of a scheduling algorithm for a network of inter-connected resources (in our case, switches) has never been studied before.

The authors of [42] propose the use of programmable hardware in the dataplane for packet scheduling and queue management, in order to achieve various objectives. The proposal shows that there is no “silver bullet” solution, by simulating three schemes (FQ, CoDel+FQ, CoDel+FIFO) competing on three different metrics. As mentioned earlier, our work is inspired by the questions the authors raise; we adopt a broader view of scheduling in which packets can carry dynamic state leading to the results presented here. A recent proposal for programmable packet scheduling [41], developed in parallel to UPS, uses an hierarchy of priority and calendar queues to express different scheduling algorithms on a single switch hardware. The proposed solution is able to achieve better expressiveness than LSTF by allowing packet headers to be re-initialized at *every switch*. UPS assumes a stronger model, where the header initialization is restricted to the ingress routers, while the core switches remain untouched. Moreover, we provide theoretical results which shed light on the effectiveness of both of these models.

7 Conclusion

This paper started with a theoretical perspective by analyzing whether there exists a single *universal* packet scheduling algorithm that can perfectly replay all viable

¹¹A scheduling algorithm is said to be optimal if it can (feasibly) schedule a set of tasks that can be scheduled by any other algorithm.

schedules. We proved that while such an algorithm cannot exist, LSTF comes closest to being one (in terms of the number of congestion points it can handle). We then empirically demonstrated the ability of LSTF to approximately replay a wide range of scheduling algorithms under varying network settings. Replaying a given schedule, while of theoretical interest, requires the knowledge of viable output times, which is not available in practice.

Hence, we next considered if LSTF can be used in practice to achieve various performance objectives. We showed via simulation how LSTF, combined with heuristics to set the slack values at the ingress, can do a reasonable job of minimizing average flow completion time, minimizing tail latencies, and achieving per-flow fairness. We also discussed some limitations of LSTF (with respect to achieving class-based fairness and traffic shaping).

Noting that scheduling is often used along with AQM to prevent queue build up, we then showed how LSTF can be used to implement versions of AQM from the network edge, with performance comparable to FQ-CoDel and to DCTCP with ECN (the state-of-the art AQM schemes for wide-area and datacenters respectively).

While an initial step towards understanding the notion of a Universal Packet Scheduling algorithm, our work leaves several theoretical questions unanswered, three of which we mention here. First, we showed existence of a UPS with omniscient header initialization, and nonexistence with limited-information initialization. *What is the least information we can use for header initialization in order to achieve universality?* Second, we showed that, in practice, the fraction of overdue packets is small, and most are only overdue by a small amount. *Are there tractable bounds on both the number of overdue packets and/or their degree of lateness?* Third, while we have a formal characterization for the scope of LSTF with respect to replaying a given schedule, and we have simulation evidence of LSTF’s ability to meet several performance objectives, we do not yet have any formal model for the scope of LSTF in meeting these objectives. *Can one describe the class of performance objectives that LSTF can meet?* Also, are there any new objectives that LSTF allows us to achieve?

8 Acknowledgments

We are thankful to Satish Rao for his helpful tips regarding the theoretical aspects of this work and to Anirudh Sivaraman for liberally sharing his insights on hardware implementation of fine-grained priorities. We would also like to thank Aisha Mushtaq, Kay Ousterhout, Aurojit Panda, Justine Sherry, Ion Stoica and our anonymous HotNets and NSDI reviewers for their thoughtful feedback. Finally, we would like to thank our shepherd Srikanth Kandula for helping shape the final version of this paper. This work was supported by Intel Research and by the National Science Foundation under Grant No. 1117161, 1343947 and 1040838.

References

- [1] Global Consortium to Construct New Cable System Linking US and Japan to Meet Increasing Bandwidth Demands. http://googlepress.blogspot.com/2008/02/global-consortium-to-construct-new_26.html.
- [2] High Capacity StrataXGStrident II Ethernet Switch Series. <http://www.broadcom.com/products/Switching/Data-Center/BCM56850-Series>.
- [3] Internet2. <http://www.internet2.edu/>.
- [4] Kathie Nichol's CoDel presented by Van Jacobson. <http://www.ietf.org/proceedings/84/slides/slides-84-tsvarea-4.pdf>.
- [5] Microsoft Invests in Subsea Cables to Connect Data-centers Globally. <http://goo.gl/GoXfxH>.
- [6] NS-2. <http://www.isi.edu/nsnam/ns/>.
- [7] NS-3. <http://www.nsnam.org/>.
- [8] Token Bucket Filters. <http://lartc.org/manpages/tc-tbf.html>.
- [9] M. Alizadeh, A. Greenberg, D. A. Maltz, J. Padhye, P. Patel, B. Prabhakar, S. Sengupta, and M. Sridharan. Data Center TCP (DCTCP). In *Proc. ACM SIGCOMM*, 2010.
- [10] M. Alizadeh, S. Yang, M. Sharif, S. Katti, N. McKeown, B. Prabhakar, and S. Shenker. pFabric: Minimal Near-optimal Datacenter Transport. In *Proc. ACM SIGCOMM*, 2013.
- [11] M. Allman. Comments on bufferbloat. *ACM SIGCOMM Computer Communication Review*, 2013.
- [12] T. Benson, A. Akella, and D. Maltz. Network Traffic Characteristics of Data Centers in the Wild. In *Proc. ACM Internet Measurement Conference (IMC)*, 2012.
- [13] R. Bhagwan and B. Lin. Fast and Scalable Priority Queue Architecture for High-Speed Network Switches. In *Proc. IEEE Infocom*, 2000.
- [14] P. Bosshart, D. Daly, G. Gibb, M. Izzard, N. McKeown, J. Rexford, C. Schlesinger, D. Talayco, A. Vahdat, G. Varghese, and D. Walker. P4: Programming Protocol-independent Packet Processors. *ACM SIGCOMM Computer Communication Review*, 2014.
- [15] P. Bosshart, G. Gibb, H.-S. Kim, G. Varghese, N. McKeown, M. Izzard, F. Mujica, and M. Horowitz. Forwarding Metamorphosis: Fast Programmable Match-action Processing in Hardware for SDN. In *Proc. ACM SIGCOMM*, 2013.
- [16] M. Casado, T. Koponen, S. Shenker, and A. Tootoonchian. Fabric: A Retrospective on Evolving SDN. In *Proc. ACM HotSDN*, 2012.
- [17] S.-T. Chuang, A. Goel, N. McKeown, and B. Prabhakar. Matching output queueing with a combined input/output-queued switch. *IEEE Journal on Selected Areas in Communications*, 1999.
- [18] B. Claise. Cisco systems NetFlow services export version 9. RFC 3954, 2004.
- [19] D. D. Clark, S. Shenker, and L. Zhang. Supporting Real-time Applications in an Integrated Services Packet Network: Architecture and Mechanism. *ACM SIGCOMM Computer Communication Review*, 1992.
- [20] A. Demers, S. Keshav, and S. Shenker. Analysis and Simulation of a Fair Queueing Algorithm. *ACM SIGCOMM Computer Communication Review*, 1989.
- [21] N. Dukkupati and N. McKeown. Why Flow-Completion Time is the Right Metric for Congestion Control. *ACM SIGCOMM Computer Communication Review*, 2006.
- [22] S. Floyd and V. Jacobson. Random Early Detection Gateways for Congestion Avoidance. *IEEE/ACM Trans. Netw.*, 1993.
- [23] G. Gibb, G. Varghese, M. Horowitz, and N. McKeown. Design principles for packet parsers. In *ACM/IEEE Symposium on Architectures for Networking and Communications Systems (ANCS)*, 2013.
- [24] A. Gupta, M. T. Hajiaghayi, and H. Räcke. Oblivious Network Design. In *Proceedings of the Seventeenth Annual ACM-SIAM Symposium on Discrete Algorithm*, SODA '06, 2006.
- [25] T. Hoeiland-Joergensen, P. McKenney, D. Taht, J. Gettys, and E. Dumazet. FlowQueue-Codel. *IETF Informational*, 2013.
- [26] A. Ioannou and M. G. H. Katevenis. Pipelined Heap (Priority Queue) Management for Advanced Scheduling in High-speed Networks. *IEEE/ACM Trans. Netw.*, 2007.
- [27] R. Jain, D.-M. Chiu, and W. Hawe. A Quantitative Measure Of Fairness And Discrimination For

- Resource Allocation In Shared Computer Systems. *CoRR*, 1998.
- [28] J. Y.-T. Leung. A new algorithm for scheduling periodic, real-time tasks. *Algorithmica*, 1989.
- [29] L. Li, D. Alderson, W. Willinger, and J. Doyle. A First-principles Approach to Understanding the Internet’s Router-level Topology. In *Proc. ACM SIGCOMM*, 2004.
- [30] C. L. Liu and J. W. Layland. Scheduling Algorithms for Multiprogramming in a Hard-Real-Time Environment. *Journal of the ACM (JACM)*, 1973.
- [31] R. Mittal, J. Sherry, S. Ratnasamy, and S. Shenker. Recursively Cautious Congestion Control. In *Proc. USENIX NSDI*, 2014.
- [32] K. Nichols and V. Jacobson. Controlling Queue Delay. *Queue*, 2012.
- [33] K. Nichols and V. Jacobson. Controlled delay active queue management: draft-nichols-tsvwg-codel-02. *Internet Requests for Comments-Work in Progress*, <http://tools.ietf.org/id/draft-nichols-tsvwg-codel-01.txt>, Tech. Rep, 2014.
- [34] 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 Trans. Netw.*, 1993.
- [35] P. Phaal, S. Panchen, and N. McKee. InMon corporation’s sFlow: A method for monitoring traffic in switched and routed networks. RFC 3176, 2001.
- [36] B. Raghavan, M. Casado, T. Koponen, S. Ratnasamy, A. Ghodsi, and S. Shenker. Software-defined Internet Architecture: Decoupling Architecture from Infrastructure. In *Proc. ACM HotNets*, 2012.
- [37] K. Ramakrishnan, S. Floyd, and D. Black. The Addition of Explicit Congestion Notification (ECN) to IP, 2001.
- [38] S. Blake and D. Black and M. Carlson and E. Davies and Z. Wang and W. Weiss. An Architecture for Differentiated Services. RFC 2475, 1998.
- [39] M. Shreedhar and G. Varghese. Efficient Fair Queueing Using Deficit Round Robin. *ACM SIGCOMM Computer Communication Review*, 1995.
- [40] A. Singh, J. Ong, A. Agarwal, G. Anderson, A. Armistead, R. Bannon, S. Boving, G. Desai, B. Felderman, P. Germano, A. Kanagala, J. Provost, J. Simmons, E. Tanda, J. Wanderer, U. Hölzle, S. Stuart, and A. Vahdat. Jupiter Rising: A Decade of Clos Topologies and Centralized Control in Google’s Datacenter Network. In *Proc. ACM SIGCOMM*, 2015.
- [41] A. Sivaraman, S. Subramanian, A. Agrawal, S. Chole, S.-T. Chuang, T. Edsall, M. Alizadeh, S. Katti, N. McKeown, and H. Balakrishnan. Towards Programmable Packet Scheduling. In *Proc. ACM HotNets*, 2015.
- [42] A. Sivaraman, K. Winstein, S. Subramanian, and H. Balakrishnan. No Silver Bullet: Extending SDN to the Data Plane. In *Proc. ACM HotNets*, 2013.
- [43] N. Spring, R. Mahajan, and D. Wetherall. Measuring ISP Topologies with Rocketfuel. In *Proc. ACM SIGCOMM*, 2002.
- [44] I. Stoica, S. Shenker, and H. Zhang. Core-stateless Fair Queueing: A Scalable Architecture to Approximate Fair Bandwidth Allocations in High-speed Networks. *IEEE/ACM Trans. Netw.*, 2003.
- [45] I. Stoica and H. Zhang. Providing Guaranteed Services Without Per Flow Management. In *Proc. ACM SIGCOMM*, 1999.
- [46] L. Zhang. Virtual Clock: A New Traffic Control Algorithm for Packet Switching Networks. *ACM SIGCOMM Computer Communication Review*, 1990.

Appendix

A Proofs for Theoretical Results

This section contains theoretical proofs for the analytical replayability results presented in §2. We begin with defining some notations used throughout in the proofs.

A.1 Notations

We use the following notations for our proofs, some of which have been already defined in the main text:

Relevant nodes:

$src(p)$: Ingress of a packet p .
 $dest(p)$: Egress of a packet p .

Relevant time notations:

$T(p, \alpha)$: Transmission time of a packet p at node α .
 $o(p, \alpha)$: Time when the first bit of p is scheduled by node α in the original schedule.
 $o(p) = o(p, dest(p)) + T(p, dest(p))$: Time when the last bit of p exits the network in the original schedule (which is non-preemptive).
 $o'(p)$: Time when the last bit of p exits the network in the replay (which may be preemptive in our theoretical arguments).
 $i(p, \alpha)$ and $i'(p, \alpha)$: Time when p arrives at node α in the original schedule and in the replay respectively.

$i(p) = i(p, \text{src}(p)) = i'(p)$: Arrival time of p at its ingress. This remains the same for both the original schedule and the replay.

$t_{\min}(p, \alpha, \beta)$: Minimum time p takes to start from node α and exit from node β in an uncongested network. It therefore includes the propagation delays and the store-and-forward delays of all links in the path from α to β and the transmission delays at α and β . Handling the edge case: $t_{\min}(p, \alpha, \alpha) = T(p, \alpha)$

$\text{slack}(p) = o(p) - i(p) - t_{\min}(p, \text{src}(p), \text{dest}(p))$: Total slack of p that gets assigned at its ingress. It denotes the amount of time p can wait in the network (excluding the time when any of its bits are getting serviced) without missing its target output time.

$\text{slack}(p, \alpha, t) = o(p) - t - t_{\min}(p, \alpha, \text{dest}(p)) + T(p, \alpha)$: Remaining slack of the last bit of p at time t when it is at node α . We derive this expression in Appendix A.4.

Other miscellaneous notations:

$\text{path}(p, \alpha, \beta)$: The ordered set of nodes and links in the path taken by p to go from α to β . The set also includes α and β as the first and the last nodes.

$\text{path}(p) = \text{path}(p, \text{src}(p), \text{dest}(p))$

$\text{pass}(\alpha)$: Set of packets that pass through node α .

A.2 Existence of a UPS under Omniscient Header Initialization

Algorithm: At the ingress, insert an n -dimensional vector in the packet header, where the i^{th} element contains $o(p, \alpha_i)$, α_i being the i^{th} hop in $\text{path}(p)$. Every time a packet p arrives at the router, the router pops the value at the head of the vector in p 's header and uses that as the priority for p (earlier values of output times get higher priority). This can perfectly replay any schedule.

Proof: We can prove that the above algorithm will result in no overdue packets (which do not meet their original schedule's target) using the following two theorems:

Theorem 1: If for any node α , $\exists p' \in \text{pass}(\alpha)$, such that using the above algorithm, the last bit of p' exits α at time $(t' > (o(p', \alpha) + T(p', \alpha)))$, then $(\exists p \in \text{pass}(\alpha) | i'(p, \alpha) \leq t'$ and $i'(p, \alpha) > o(p, \alpha))$.

Proof by contradiction: Consider the first such $p^* \in \text{pass}(\alpha)$ that gets late at α (i.e. its last bit exits α at time $t^* > (o(p^*, \alpha) + T(p^*, \alpha))$). Suppose the above condition is not true i.e. $(\forall p \in \text{pass}(\alpha) | i'(p, \alpha) \leq o(p, \alpha)$ or $i'(p, \alpha) > t^*)$. In other words, if p arrives at or before time t^* , it also arrives at or before time $o(p, \alpha)$. Given that all bits of p^* arrive at or before time t^* , they also arrive at or before time $o(p^*, \alpha)$. The only reason why the last bit of p^* would wait until time $(t^* > o(p^*, \alpha) + T(p^*, \alpha))$ in our work-conserving replay is if some other bits (belonging to higher priority packets) were being scheduled after time $o(p^*, \alpha)$, resulting in p^* not being able to complete its transmission by time $(o(p^*, \alpha) + T(p^*, \alpha))$. However, as per our algo-

rithm, any packet p_{high} having a higher priority than p^* at α must have been scheduled before p^* in the original schedule, implying that $(o(p_{\text{high}}, \alpha) + T(p_{\text{high}}, \alpha)) \leq o(p^*, \alpha)$.¹² Therefore, some bits of p_{high} being scheduled after time $o(p^*, \alpha)$, implies them being scheduled after time $(o(p_{\text{high}}, \alpha) + T(p_{\text{high}}, \alpha))$. This means that p_{high} is already late and contradicts our assumption that p^* is the first packet to get late. Hence, Theorem 1 is proved by contradiction.

Theorem 2: $\forall \alpha, (\forall p \in \text{pass}(\alpha) | i'(p, \alpha) \leq i(p, \alpha))$.

Proof by contradiction: Consider the first time when some packet p^* arrives late at some node α^* (i.e. $i'(p^*, \alpha^*) > i(p^*, \alpha^*)$). In other words, α^* is the first node in the network to see a late packet arrival, and p^* is the first late arriving packet. Let α_{prev} be the node visited by p^* just before arriving at α^* . p^* can arrive at a time later than $i(p^*, \alpha^*)$ at α^* only if the last bit of p^* exits α_{prev} at time $t_{\text{prev}} > o(p^*, \alpha_{\text{prev}}) + T(p^*, \alpha_{\text{prev}})$. As per Theorem 1 above, this is possible only if some packet p' (which may or may not be the same as p^*) arrives at α_{prev} at time $i'(p', \alpha_{\text{prev}}) > o(p', \alpha_{\text{prev}}) \geq i(p', \alpha_{\text{prev}})$ and $i'(p', \alpha_{\text{prev}}) \leq t_{\text{prev}} < i'(p^*, \alpha^*)$. This contradicts our assumption that α^* is the first node to see a late arriving packet. Therefore, $\forall \alpha, (\forall p \in \text{pass}(\alpha) | i'(p, \alpha) \leq i(p, \alpha))$.

Combining the two theorems above: Since $\forall \alpha (\forall p \in \text{pass}(\alpha) | i'(p, \alpha) \leq i(p, \alpha))$, with the above algorithm, $\forall \alpha (\forall p \in \text{pass}(\alpha))$, all bits of p exit α before $(o(p, \alpha) + T(p, \alpha))$. Therefore, the algorithm can perfectly replay any viable schedule.

A.3 Nonexistence of a UPS under black-box initialization

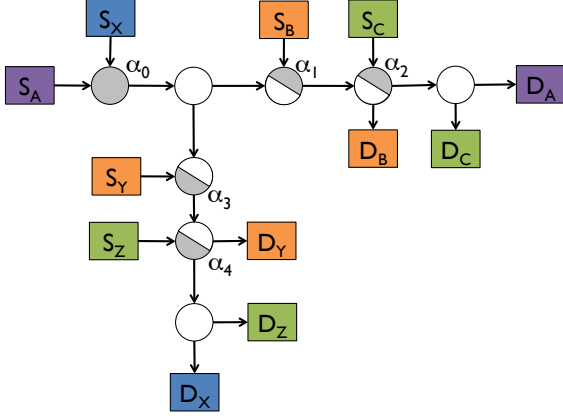
Proof by counter-example: Consider the example shown in Figure 7. For simplicity, assume all the propagation delays are zero, the transmission time for each congestion point (shaded in gray) is 1 unit and the uncongested (white) routers have zero transmission time.¹³ All packets are of the same size.

The table illustrates two cases. For each case, a packet's arrival and scheduling time (the time when the packet is scheduled by the router) at each node through which it passes are listed. A packet represented by p belongs to flow P , with ingress S_P and egress D_P , where $P \in \{A, B, C, X, Y, Z\}$. The packets have the same path in both cases. For example, a belongs to Flow A, starts at ingress S_A , exits at egress D_A and passes through three congestion points in its path α_0 , α_1 and α_2 ; x belongs to Flow X, starts at ingress S_X , exits at egress D_X and passes through three congestion points in its path α_0 , α_3 and α_4 ; and so on.

The two critical packets we care about in this example

¹²Given that the original schedule is non-preemptible, the next packet gets scheduled only after the previous one has completed its transmission.

¹³These assignments are made for simplicity of understanding. The example will hold for any reasonable value of propagation and transmission delays.



Node	Packet(arrival time, scheduling time)
<i>Case 1</i>	
α_0	$a(\mathbf{0},0); x(\mathbf{0},1)$
α_1	$a(1,1), b_1(2,2), b_2(3,3), b_3(4,4)$
α_2	$c_1(2,2), c_2(3,3); a(2,\mathbf{4})$
α_3	$x(2,2), y_1(2,3), y_2(3,4)$
α_4	$z(2,2), x(3,\mathbf{3})$
<i>Case 2</i>	
α_0	$x(\mathbf{0},0); a(\mathbf{0},1)$
α_1	$a(2,2), b_1(2,3), b_2(3,4), b_3(4,5)$
α_2	$c_1(2,2), c_2(3,3), a(3,\mathbf{4})$
α_3	$x(1,1), y_1(2,2), y_2(3,3)$
α_4	$z(2,2), x(2,\mathbf{3})$

Figure 7: Example showing non-existence of a UPS with Blackbox Initialization. A packet represented by p belongs to flow P , with ingress S_P and egress D_P , where $P \in \{A, B, C, X, Y, Z\}$. For simplicity assume all packets are of the same size and all links have a propagation delay of zero. All uncongested routers (white), ingresses and egresses have a transmission time of zero. The congestion points (shaded gray) have transmission times of $T = 1$ unit.

are a and x , which interact with each-other at their first congestion point α_0 , being scheduled by α_0 at different times in the two cases (a before x in Case 1 and x before a in Case 2). But, notice that for both cases,

1. a enters the network from its ingress S_A at congestion point α_0 at time 0, and passes through two other congestion points α_1 and α_2 before exiting the network at time $(4+1)$ ¹⁴.
2. x enters the network from its ingress S_X at congestion point α_0 at time 0, and passes through two other congestion points α_3 and α_4 before exiting the network at time $(3+1)$.

a interacts with packets from Flow C at its third congestion point α_2 , while x interacts with a packet from Flow Z at its third congestion point α_4 . For both cases,

1. Two packets of Flow C (c_1, c_2) enter the network at times 2 and 3 at α_2 before they exit the network at time $(2+1)$ and $(3+1)$ respectively.
2. z enters the network at time 2 at α_4 before exiting at

¹⁴ $+1$ is added to indicate transmission time at the last congestion point. As mentioned before, we assume the propagation delay to the egress and the transmission time at the egress are both 0.

time $2+1$.

The difference between the two cases comes from how a interacts with packets from Flow B at its second congestion point α_1 and how x interacts with packets from Flow Y at its second congestion points α_3 . Note that α_1 and α_3 are the last congestion points for Flow B and Flow Y packets respectively and their exit times from these congestion points directly determine their exit times from the network.

1. Three packets of Flow B (b_1, b_2, b_3) enter the network at times 2, 3 and 4 respectively at α_1 . In Case 1, they leave α_1 at times $(2+1), (3+1), (4+1)$ respectively. This provides no *lee-way* for a at α_0 , which leaves α_1 at time $(1+1)$, since it is required that α_1 must schedule a by at most time 3 in order for it to exit the network at its target output time. In Case 2, (b_1, b_2, b_3) leave at times $(3+1), (4+1), (5+1)$ respectively, providing *lee-way* for a at α_0 , which leaves α_1 at time $(2+1)$.
2. Two packets of Flow Y (y_1, y_2) enter the network at times 2 and 3 respectively at α_3 . In Case 1, they leave at times $(3+1), (4+1)$ respectively, providing a *lee-way* for x at α_0 , which leaves α_3 at time $(2+1)$. In Case 2, (y_1, y_2) exit at times $(2+1), (3+1)$, providing no *lee-way* for x at α_0 , which leaves α_3 at time $(1+1)$.

Note that the interaction of a and x with Flow C and Flow Z at their third congestion points respectively, is what ensures that their eventual exit time remains the same across the two cases inspite of the differences in how a and x are scheduled in their previous two hops.

Thus, we can see that $i(a), o(a), i(x), o(x)$ are the same in both cases (also indicated in bold blue). Yet, *Case 1* requires a to be scheduled before x at α_0 , else packets will get delayed at α_1 , since it is required that α_1 schedules a at a time no more than 3 units if it is to meet its target output time. *Case 2* requires x to be scheduled before a at α_0 , else packets will be delayed at α_3 , where it is required to schedule x at a time no more than 2 units if it is to meet its target output time. Since the attributes $(i(\cdot), o(\cdot), path(\cdot))$ for both a and x are exactly the same in both cases, any deterministic UPS with Blackbox Initialization will produce the same order for the two packets at α_0 , which contradicts the situation where we want a before x in one case and x before a in another.

A.4 Deriving the Slack Equation

We now prove that for any packet p waiting at any node α at time t_{now} , the remaining slack of the last bit of p is given by $slack(p, \alpha, t_{now}) = o(p) - t_{now} - t_{min}(p, \alpha, dest(p)) + T(p, \alpha)$.

Let $t_{wait}(p, \alpha, t_{now})$ denote the total time spent by p on waiting behind other packets at the nodes in its path from $src(p)$ to α (including these two nodes) until time t_{now} . We define $t_{wait}(p, \alpha, t_{now})$, such that it excludes the transmission times at previous nodes which gets captured in t_{min} , but includes the local service time received by the

packet so far at α itself.

$$\text{slack}(p, \alpha, t_{\text{now}}) = \text{slack}(p) - t_{\text{wait}}(p, \alpha, t_{\text{now}}) + T(p, \alpha) \quad (1a)$$

$$= o(p) - i(p) - t_{\min}(p, \text{src}(p), \text{dest}(p)) - t_{\text{wait}}(p, \alpha, t_{\text{now}}) + T(p, \alpha) \quad (1b)$$

$$= o(p) - i(p) - (t_{\min}(p, \text{src}(p), \alpha) + t_{\min}(p, \alpha, \text{dest}(p)) - T(p, \alpha)) - t_{\text{wait}}(p, \alpha, t_{\text{now}}) + T(p, \alpha) \quad (1c)$$

$$= o(p) - t_{\min}(p, \alpha, \text{dest}(p)) + T(p, \alpha) - (i(p) + t_{\min}(p, \text{src}(p), \alpha) - T(p, \alpha) + t_{\text{wait}}(p, \alpha, t_{\text{now}})) \quad (1d)$$

$$= o(p) - t_{\min}(p, \alpha, \text{dest}(p)) + T(p, \alpha) - t_{\text{now}} \quad (1e)$$

(1a) is straightforward from our definition of LSTF and how the slack gets updated at every time slice. $T(p, \alpha)$ is added since α needs to locally consider the slack of the last bit of the packet in a store-and-forward network. (1c) then uses the fact that for any α in $\text{path}(p)$, $(t_{\min}(p, \text{src}(p), \text{dest}(p)) = t_{\min}(p, \text{src}(p), \alpha) + t_{\min}(p, \alpha, \text{dest}(p)) - T(p, \alpha))$. $T(p, \alpha)$ is subtracted here as it is accounted for twice when we break up the equation for $t_{\min}(p, \text{src}(p), \text{dest}(p))$. (1e) then follows from the fact that the difference between t_{now} and $i(p)$ is equal to the total amount of time the packet has spent in the network until time t_{now} i.e. $(t_{\text{now}} - i(p) = (t_{\min}(p, \text{src}(p), \alpha) - T(p, \alpha)) + t_{\text{wait}}(p, \alpha, t_{\text{now}}))$. We need to subtract $T(p, \alpha)$, since by our definition, $t_{\min}(p, \text{src}(p), \alpha)$ includes transmission time of the packet at α .

A.5 LSTF and EDF Equivalence

In our network-wide extension of EDF scheduling, every router computes a deadline (or priority) for a packet p based on the static header value $o(p)$ and additional state information about the minimum time the packet would take to reach its destination from the router. More precisely, each router (say α), uses $\text{priority}(p) = (o(p) - t_{\min}(p, \alpha, \text{dest}(p)) + T(p, \alpha))$ to do priority scheduling, with $o(p)$ being the value carried by the packet header, initialized at the ingress and remaining unchanged throughout. EDF is equivalent to LSTF, in that for a given original schedule, the two produce exactly the same replay schedule.

Proof: Consider a node α and let $P(\alpha, t_{\text{now}})$ be the set of packets waiting at the output queue of α at time t_{now} . A packet will then be scheduled by α as follows:

With EDF: Schedule packet $p_{\text{edf}}(\alpha, t_{\text{now}})$, where

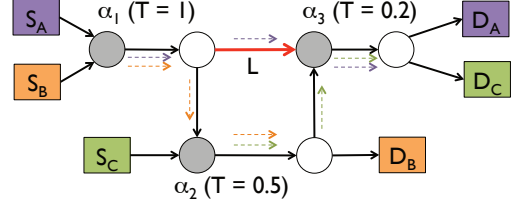
$$p_{\text{edf}}(\alpha, t_{\text{now}}) = \underset{p \in P(\alpha, t_{\text{now}})}{\text{argmin}} (\text{priority}(p, \alpha))$$

$$\text{priority}(p, \alpha) = o(p) - t_{\min}(p, \alpha, \text{dest}(p)) + T(p, \alpha)$$

With LSTF: Schedule packet $p_{\text{lsth}}(\alpha, t_{\text{now}})$, where

$$p_{\text{lsth}}(\alpha, t_{\text{now}}) = \underset{p \in P(\alpha, t_{\text{now}})}{\text{argmin}} (\text{slack}(p, \alpha, t_{\text{now}}))$$

$$\text{slack}(p, \alpha, t_{\text{now}}) = o(p) - t_{\min}(p, \alpha, \text{dest}(p)) + T(p, \alpha) - t_{\text{now}}$$



Node	Packet(arrival time, scheduling time)
α_1	$a(0,0), b(0,1)$
α_2	$b(2,2), c(2,2.5)$
α_3	$c(3,3), a(3,3.2)$

Figure 8: Example showing replay failure with simple priorities for a schedule with two congestion points per packet. A packet represented by p belongs to flow P , with ingress S_P and egress D_P , where $P \in \{A, B, C\}$. All packets are of the same size. For simplicity assume all links (except L) have a propagation delay of zero. L has a propagation delay of 2. All uncongested routers (white circles), ingresses and egresses have a transmission time of zero. The three congestion points – $\alpha_1, \alpha_2, \alpha_3$ have transmission times of $T = 1$ unit, $T = 0.5$ units and $T = 0.2$ units respectively.

The above expression for $\text{slack}(p, \alpha, t_{\text{now}})$ has been derived in §A.4. Thus, $\text{slack}(p, \alpha, t_{\text{now}}) = \text{priority}(p, \alpha) - t_{\text{now}}$. Since t_{now} is the same for all packets, we can conclude that:

$$\underset{p \in P(\alpha, t_{\text{now}})}{\text{argmin}} (\text{slack}(p, \alpha, t_{\text{now}})) = \underset{p \in P(\alpha, t_{\text{now}})}{\text{argmin}} (\text{priority}(p, \alpha))$$

$$\implies p_{\text{lsth}}(\alpha, t_{\text{now}}) = p_{\text{edf}}(\alpha, t_{\text{now}})$$

Therefore, at any given point of time, all nodes will schedule the same packet with both EDF and LSTF (assuming ties are broken in the same way for both EDF and LSTF, such as by using FCFS). Hence, EDF and LSTF are equivalent.

A.6 Simple Priorities Replay Failure for Two Congestion Points Per Packet

In Figure 8, we present an example which shows that simple priorities can fail in replay when there are two congestion points per packet, no matter what information is used to assign priorities. At α_1 , we need to have $\text{priority}(a) < \text{priority}(b)$, at α_2 we need to have $\text{priority}(b) < \text{priority}(c)$ and at α_3 we need to have $\text{priority}(c) < \text{priority}(a)$. This creates a priority cycle where we need $\text{priority}(a) < \text{priority}(b) < \text{priority}(c) < \text{priority}(a)$, which can never be possible to achieve with simple priorities.

We would also like to point out here that priority assignment for perfect replay in networks with single congestion point per packet requires detailed knowledge about the topology and the input load. More precisely, if a packet p passes through congestion point α_p , then its priority needs to be assigned as $\text{priority}(p) = o(p) - t_{\min}(p, \alpha_p, \text{dest}(p)) + T(p, \alpha_p)$. The proof that this would always replay schedules with at most one congestion point per packet follows from the fact that the only scheduling decision made in a packet p 's path is at the single conges-

tion point α_p . This decision, at the single congestion point in a packet's path, is the same as what will be made with the network-wide extension of EDF, which we proved is equivalent to LSTF in §A.5. LSTF, in turn, can always replay schedules with one (or to be more precise, at most two) congestion points per packet, as we shall prove in §A.7.

Hence, in order to replay schedules with at most one congestion per packet using simple priorities, we need to know where the congestion point occurs in a packet's path, along with the final output times, to assign the priorities. In the absence of this knowledge, priorities cannot replay even a single congestion point.

A.7 LSTF: Perfect Replay for at most Two Congestion Points per Packet

A.7.1 Main Proof

We now prove that LSTF can replay all schedules with at most two congestion points per packet. Note that we work with bits in our proof, since we assume a preemptive version of LSTF. Due to store-and-forward routers, the remaining slack of a packet at a particular router is represented by the slack of the last bit of the packet (with all other bits of the packet having the same slack as the last bit).

In order for a replay failure to occur, there must be at least one overdue packet, where a packet p is said to be overdue if $o'(p) > o(p)$. This implies that p must have spent all of its slack while waiting behind other packets at a queue in some node α at say time t , such that $slack(p, \alpha, t) < 0$. Obviously, α must be a congestion point.

Necessary Condition for Replay Failure with LSTF:

If a packet p^* sees negative slack at a congestion point α when its last bit exits α at time t^* in the replay (i.e. $slack(p^*, \alpha, t^*) < 0$), then $(\exists p \in pass(\alpha) \mid i'(p, \alpha) \leq t^*$ and $i'(p, \alpha) > o(p, \alpha))$. We prove this in §A.7.2.

We use the term “local deadline of p at α ” for $o(p, \alpha)$, which is the time at which α schedules p in the original schedule.

Key Observation: *When there are at most two congestion points per packet, then no packet p can arrive at any congestion point α in the replay, after its local deadline at α (i.e. $i'(p, \alpha) > o(p, \alpha)$ is not possible). Therefore, by the necessary condition above, no packet can see a negative slack at any congestion point.*

Proof by contradiction: Suppose that there exists α^* , which is the first congestion point (in time) that sees a packet which arrives after its local deadline at α^* . Let p^* be this first packet that arrives after its local deadline at α^* ($i'(p^*, \alpha^*) > o(p^*, \alpha^*)$). Since there are at most two congestion points per packet, either α^* is the first congestion point seen by p^* or the last (or both).

(1) If α^* is the first congestion point seen by p^* , then clearly, $i'(p^*, \alpha^*) = i(p^*, \alpha^*) \leq o(p^*, \alpha^*)$. This contradicts our assumption that $i'(p^*, \alpha^*) > o(p^*, \alpha^*)$.

(2) If α^* is not the first congestion point seen by p^* , then it is the last congestion point seen by p^* . If $i'(p^*, \alpha^*) > o(p^*, \alpha^*)$, then it would imply that p^* saw a negative slack before arriving at α^* . Suppose p^* saw a negative slack at a congestion point α_{prev} , before arriving at α^* when its last bit exited α_{prev} at time t_{prev} . Clearly, $t_{prev} < i'(p^*, \alpha^*)$. As per our necessary condition, this would imply that there must be another packet p' , such that $i'(p', \alpha_{prev}) > o(p', \alpha_{prev})$ and $i'(p', \alpha_{prev}) \leq t_{prev} < i'(p^*, \alpha^*)$. This contradicts our assumption that α^* is the first congestion point (in time) that sees a packet which arrives after its corresponding scheduling time in the original schedule.

Hence, no congestion point can see a packet that arrives after its local deadline at that congestion point (and therefore no packet can get overdue) when there are at most two congestion points per packet.

A.7.2 Proof for Necessary Condition for Replay Failure with LSTF

We start this proof with the following observation:

Observation 1: If all bits of a packet p exit a router α by time $o(p, \alpha) + T(p, \alpha)$, then p cannot see a negative slack at α .

Proof for Observation 1: As shown previously in §A.4,

$$slack(p, \alpha, t) = o(p) - t_{min}(p, \alpha, dest(p)) + T(p, \alpha) - t$$

Therefore,

$$\begin{aligned} & slack(p, \alpha, o(p, \alpha) + T(p, \alpha)) \\ &= o(p) - t_{min}(p, \alpha, dest(p)) + T(p, \alpha) - (o(p, \alpha) + T(p, \alpha)) \\ \text{But, } & o(p) = o(p, \alpha) + t_{min}(p, \alpha, dest(p)) + wait(p, \alpha, dest(p)) \\ \implies & slack(p, \alpha, o(p, \alpha) + T(p, \alpha)) = wait(p, \alpha, dest(p)) \\ \implies & slack(p, \alpha, o(p, \alpha) + T(p, \alpha)) \geq 0 \end{aligned}$$

where $wait(p, \alpha, dest(p))$ is the time spent by p in waiting behind other packets in the original schedule, after it left α , which is clearly non-negative.

We now move to the main proof for the necessary condition.

Necessary Condition for Replay Failure: If a packet p^* sees negative slack at a congestion point α when its last bit exits α at time t^* in the replay (i.e. $slack(p^*, \alpha, t^*) < 0$), then $(\exists p \in pass(\alpha) \mid i'(p, \alpha) \leq t^*$ and $i'(p, \alpha) > o(p, \alpha))$.

Proof by Contradiction: Suppose this is not the case i.e. there exists p^* whose last bit exits α at time t^* , such that $slack(p^*, \alpha, t^*) < 0$ and $(\forall p \in pass(\alpha) \mid i'(p, \alpha) > t^*$ or $i'(p, \alpha) \leq o(p, \alpha))$. We can show that if the latter condition holds, then p^* cannot see a negative slack at α , thus violating our assumption.

We take the set of all bits which exit α at or before time t^* in the LSTF replay schedule. We denote this set as $S_{bits}(\alpha, t^*)$. As per our assumption, $(\forall b \in S_{bits}(\alpha, t^*) \mid i'(p_b, \alpha) \leq o(p_b, \alpha))$, where p_b denotes the packet to which

bit b belongs. Note that $S_{bits}(\alpha, t^*)$ also includes all bits of p^* , since they all arrive before time t^* .

We now prove that no bit in $S_{bits}(\alpha, t^*)$ can see a negative slack (and therefore p^* cannot see a negative slack at α), leading to a contradiction. The proof comprises of two steps:

Step 1: Using the same input arrival times of each packet at α as in the replay schedule, we first construct a *feasible schedule* at α up until time t^* , denoted by $FS(\alpha, t^*)$, where by feasibility we mean that no bit in $S_{bits}(\alpha, t^*)$ sees a negative slack.

Step 2: We then do an iterative transformation of $FS(\alpha, t^*)$ such that the bits in $S_{bits}(\alpha, t^*)$ are scheduled in the order of their *least remaining slack times*. This reproduces the LSTF replay schedule from which $FS(\alpha, t^*)$ was constructed in the first place. However, while doing the transformation we show how the schedule remains feasible at every iteration, proving that the LSTF schedule finally obtained is also feasible up until time t^* . In other words, no packet sees a negative slack at α in the resulting LSTF replay schedule up until time t^* , contradicting our assumption that p^* sees a negative slack when it exits α at time t^* in the replay.

We now discuss these two steps in details.

Step 1: Construct a feasible schedule at α up until time t^* (denoted as $FS(\alpha, t^*)$) for which no bit in $S_{bits}(\alpha, t^*)$ sees a negative slack.

(i) Algorithm for constructing $FS(\alpha, t^*)$: Use priorities to schedule each bit in $S_{bits}(\alpha, t^*)$, where $\forall b \in S_{bits}(\alpha, t^*) \mid priority(b) = o(p_b, \alpha)$. (Note that since both $FS(\alpha, t^*)$ and LSTF are work-conserving, $FS(\alpha, t^*)$ is just a shuffle of the LSTF schedule up until t^* . The set of time slices at which a bit is scheduled in $FS(\alpha, t^*)$ and in the LSTF schedule up until t^* remains the same, but *which* bit gets scheduled at a given time slice is different.)

(ii) In $FS(\alpha, t^*)$, all bits b in $S_{bits}(\alpha, t^*)$ exit α by time $o(p_b, \alpha) + T(p_b, \alpha)$.

Proof by contradiction: Suppose the statement is not true and consider the first bit b^* that exits after time $(o(p_{b^*}, \alpha) + T(p_{b^*}, \alpha))$. We term this as b^* got late at α due to $FS(\alpha, t^*)$. Remember that, as per our assumption, $(\forall b \in S_{bits}(\alpha, t^*) \mid i'(p_b, \alpha) \leq o(p_b, \alpha))$. Thus, given that all bits of p_{b^*} arrive at or before time $o(p_{b^*}, \alpha)$, the only reason why the delay can happen in our work-conserving $FS(\alpha, t^*)$ is if some other higher priority bits were being scheduled after time $o(p_{b^*}, \alpha)$, resulting in p_{b^*} not being able to complete its transmission by time $(o(p_{b^*}, \alpha) + T(p_{b^*}, \alpha))$. However, as per our priority assignment algorithm, any bit b' having a higher priority than b^* at α must have been scheduled before the first bit of p_{b^*} in the non-preemptible original schedule, implying that $(o(p_{b'}, \alpha) + T(p_{b'}, \alpha)) \leq o(p_{b^*}, \alpha)$. Therefore, a bit b' being scheduled after time $o(p_{b^*}, \alpha)$, implies it being scheduled after time $(o(p_{b'}, \alpha) + T(p_{b'}, \alpha))$. This contradicts our assumption that b^* is the first bit to get late at α due to $FS(\alpha, t^*)$. Therefore, all bits

b in $S_{bits}(\alpha, t^*)$ exit α by time $o(p_b, \alpha) + T(p_b, \alpha)$ as per the schedule $FS(\alpha, t^*)$.

(iii) Since all bits in $S_{bits}(\alpha, t^*)$ exit by time $o(p_b, \alpha) + T(p_b, \alpha)$ due to $FS(\alpha, t^*)$, no bit in $S_{bits}(\alpha, t^*)$ sees a negative slack at α (from Observation 1).

Step 2: Transform $FS(\alpha, t^*)$ into a feasible LSTF schedule for the single switch α up until time t^* .

(Note: The following proof is inspired from the standard LSTF optimality proof that shows that for a single switch, any feasible schedule can be transformed to an LSTF (or EDF) schedule [30]).

Let $fs(b, \alpha, t^*)$ be the scheduling time slice for bit b in $FS(\alpha, t^*)$. The transformation to LSTF is carried out by the following pseudo-code:

```

1: while true do
2:   Find two bits,  $b_1$  and  $b_2$ , such that:
      ( $fs(b_1, \alpha, t^*) < fs(b_2, \alpha, t^*)$ ) and
      ( $slack(b_2, \alpha, fs(b_1, \alpha, t^*))$ 
        $< slack(b_1, \alpha, fs(b_1, \alpha, t^*))$ ) and
      ( $i'(b_2, \alpha, t^*) \leq fs(b_1, \alpha, t^*)$ )
3:   if no such  $b_1$  and  $b_2$  exist then
4:      $FS(\alpha, t^*)$  is an LSTF schedule
5:     break
6:   else
7:      $swap(fs(b_1, \alpha, t^*), fs(b_2, \alpha, t^*))$  ▷
      swap the scheduling times of the two bits. 15
8:   end if
9: end while
10: Shuffle the scheduling time of the bits belonging to
     the same packet, to ensure that they are in order.
11: Shuffle the scheduling time of the same-slack bits
     such that they are in FIFO order

```

Line 7 above will not cause b_1 to have a negative slack, when it gets scheduled at $fs(b_2, \alpha, t^*)$ instead of $fs(b_1, \alpha, t^*)$. This is because the difference in $slack(b_2, \alpha, t)$ and $slack(b_1, \alpha, t)$ is independent of t and so:

$$\begin{aligned} & slack(b_2, \alpha, fs(b_1, \alpha, t^*)) < slack(b_1, \alpha, fs(b_1, \alpha, t^*)) \\ \implies & slack(b_2, \alpha, fs(b_2, \alpha, t^*)) < slack(b_1, \alpha, fs(b_2, \alpha, t^*)) \end{aligned}$$

Since $FS(\alpha, t^*)$ is feasible before the swap, $slack(b_2, \alpha, fs(b_2, \alpha, t^*)) \geq 0$. Therefore, $slack(b_1, \alpha, fs(b_2, \alpha, t^*)) > 0$ and the resulting $FS(\alpha, t^*)$ after the swap remains feasible.

Lines 10 and 11 will also not result in any bit getting a negative slack, because all bits participating in the shuffle have the same slack at any fixed point of time in α .

Therefore, no bit in $S_{bits}(\alpha, t^*)$ has a negative slack at α after any iteration.

Since no bit in $S_{bits}(\alpha, t^*)$ has a negative slack at α in the swapped LSTF schedule, it contradicts our statement

¹⁵Note that we are working with bits here for easy expressibility. In practice, such a swap is possible under the preemptive LSTF model.

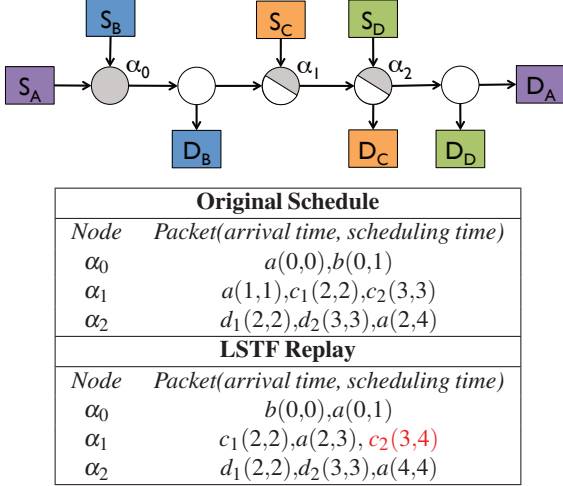


Figure 9: Example showing replay failure with LSTF when there is a flow with three congestion points. A packet represented by p belongs to flow P , with ingress S_P and egress D_P , where $P \in \{A, B, C, D\}$. For simplicity assume all links have a propagation delay of zero. All uncongested routers (white), ingresses and egresses have a transmission time of zero. The three congestion points (shaded gray) have transmission times of $T = 1$ unit

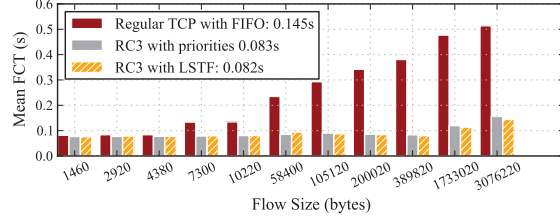
that p^* sees a negative slack when its last bit exits α at time t^* . Hence proved that if a packet p^* sees a negative slack at congestion point α when its last bit exits α at time t^* in the replay, then there must be at least one packet that arrives at α in the replay at or before time t^* and later than the time at which it is scheduled by α in the original schedule.

A.7.3 Replay Failure Example with LSTF

In Figure 9, we present an example where a flow passes through three congestion points and a replay failure occurs with LSTF. When packet a arrives at α_0 , it has a slack of 2 (since it waits behind d_1 and d_2 at α_2), while at the same time, packet b has a slack of 1 (since it waits behind a at α_0). As a result, b gets scheduled before a in the LSTF replay. a therefore arrives at α_1 with slack 1 at time 2. c_1 with a zero slack is prioritized over a . This reduces a 's slack to zero at time 3, when c_2 is also present at α_1 with zero slack. Scheduling a before c_2 , will result in c_2 being overdue (as shown). Likewise, scheduling c_2 before a would have resulted in a getting overdue. Note that in this failure case, a arrives at α_1 at time 2, which is greater than $o(a, \alpha_1) = 1$.

B Minimizing Average FCT by using RC3 with LSTF

We now look at how in-network scheduling can be used along with changes in the endhost TCP stack to minimize average flow completion times. We use RC3 [31] as our comparison-point for this objective (as it has better performance than RCP [21] and is simple to implement). In RC3 the senders aggressively send additional packets to quickly use up the available network capacity, but these packets are sent at lower priority levels to ensure that the regular traffic is not penalized. Therefore, it allows



Expt. Setup	Avg FCT (s)		
	TCP-FIFO	RC3-priorities	RC3-LSTF
12 1Gbps-10Gbps at 30% util.	0.145	0.083	0.082
12 1Gbps-10Gbps at 50% util.	0.159	0.094	0.089
12 1Gbps-10Gbps at 70% util.	0.180	0.107	0.102
12 1Gbps-1Gbps at 30% util.	0.134	0.075	0.073
12 / 10 at 30% util.	0.32	0.215	0.233
Rocketfuel at 30% util.	0.171	0.102	0.101

Figure 10: The graph shows the mean FCT bucketed by flow size for the 12 1Gbps-10Gbps topology with 30% utilization for regular TCP using FIFO and for RC3 using priorities and LSTF. The legend indicates the mean FCT across all flows. The table indicates the mean FCTs for varying settings.

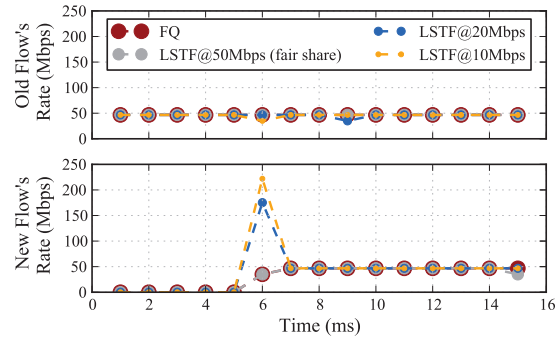


Figure 11: 20 flows share a single bottleneck link of 1Gbps and a 21st flow is added after 5ms. The graph shows the rate allocations for an old flow and the new flow with Fair Queuing and for LSTF with varying r_{est} .

near-optimal bandwidth utilization, while maintaining the cautiousness of TCP.

Slack Initialization: The slack for a packet p is initialized as $slack(p) = prio_{rc3} * D$, where $prio_{rc3}$ is the priority of the packet assigned by RC3 and D is a value much larger than the queuing delay seen by any packet in the network. We use a value of $D = 1$ sec for our simulations.

Evaluation: To evaluate RC3 with LSTF, we reuse the ns-3 [7] implementation of RC3 (along with the same TCP parameters used by RC3, such as an initial congestion window of 4), and implement LSTF in ns-3. Figure 10 shows our results. We see that using LSTF with RC3 performs comparable to (and often slightly better than) using priorities with RC3, both giving significantly lower FCTs than regular TCP with FIFO.

C Fairness Deep Dive

C.1 Understanding how LSTF provides long-term fairness

The reason behind why any slack assignment with $r_{est} < r^*$ leads to convergence to fairness is quite straight-forward and is explained by the control experiment shown in

Figure 11. 20 long-lived TCP flows share a single bottleneck link of 1Gbps (giving a fair share rate of 50Mbps) and a 21st flow is added after 5ms. Since the first 20 flows have started early, the queue at the bottleneck link already contains packets belonging to these flows.

When $r_{est} = 50Mbps$, the actual queuing delay experienced by a packet is almost equal to the slack value assigned to it. Therefore, at any given point of time, the first packet of each flow present in the queue will have a slack value which is approximately equal to zero. The next packet of each flow will have a higher slack value (around $1500bytes/50Mbps = 0.24ms$). By the time the corresponding first packets of every flow in the queue have been transmitted, the slack values of the next packet would also have been reduced to zero and so on. It therefore produces a round-robin pattern for scheduling packets across flows, as is done by FQ. Therefore, when the 21st flow starts at 5ms, with the first packet coming in with zero slack, the next one with 0.24ms slack and so on, it immediately starts following the round-robin pattern as well.

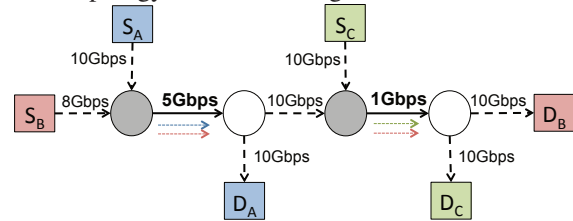
However, when r_{est} is smaller than 50Mbps, then the packets of the old flows already present in the queue have a higher slack value than what they actually experience in the network. The first packet of every flow in the queue therefore has a slack which is more than 0 when the 21st flow comes in at 5ms. The earlier packets of the new flow therefore get precedence over any of the existing packets of the old flows, resulting in the spike in the rate allocated to the new flow as shown in Figure 11. Nonetheless, with the slack of every newly arriving packet of the 21st flow being higher than the previous one and with the slack of the already queued up packet decreasing with time, the slack value of the first packet in the queue for new flow and the old flows soon *catch up* with each other and the schedule starts following a round robin pattern again. The closer r_{est} is to the fair-share rate, the sooner the slack values of the old flows and the new flow *catch up* with each other. The duration for which a packet ends up waiting in the queue is upper-bounded by the time it would have waited, had all the flows arrived at the same time and were being serviced at their fair share rate.

C.2 Weighted Fairness with multiple-bottlenecks

One can see how the above logic can be extended for achieving weighted fairness. Moreover, when a packet sees multiple bottlenecks, the slack update (subtraction of the duration for which the packet waits) at the first bottleneck ensures that the next bottleneck takes into account the rate-limiting happening at the first one and the packets are given precedence accordingly.

We did a control experiment to evaluate weighted fairness with LSTF on a multi-bottleneck topology. We started three UDP flows with a start-time jitter between 0 and 1ms,

on the topology as shown in Figure 12. We ran the sim-



r_{est} value (Mbps)			Expected Throughput (Mbps)			Achieved Throughput (Mbps)		
A	B	C	A	B	C	A	B	C
2000	100	100	4761	238	762	4762	238	763
900	100	100	4500	500	500	4499	501	500
500	100	100	4167	500	500	4166	501	500
200	100	100	3333	500	500	3333	501	500
100	100	100	2500	500	500	2500	500	501
100	100	500	2500	167	833	2500	167	834

Figure 12: Weighted Fairness on a multi-bottleneck topology (drawn above). The link capacities and the source/destination of each flow are indicated in the figure. Flows A and B share a 5Gbps link and then Flows B and C share a 1Gbps link.

Refresh Threshold (ms)	Avg FCT across bytes (s)	Avg RTT across bytes (s)
10	3.578	0.143
20	3.739	0.139
30	3.954	0.135
40	4.079	0.132

Table 4: Effect of varying refresh threshold on I2/10 topology at 70% utilization running LSTF ($r_{est} = 10Mbps$) with Edge-CoDel.

ulation for 30ms and computed the throughput each flow received for the last 15ms. We varied the values of r_{est} used for assigning slacks to each flow, relative to one another, to assign different weights to different flows. For example, r_{est} assignment $\{A: 900Mbps, B: 100Mbps, C: 100Mbps\}$ results in Flow A getting 9 times more share on the 5Gbps link than Flow B, with Flows B and C sharing the 1Gbps link equally. We compute the expected throughput based on the assigned r_{est} values and find that the throughput actually achieved is almost the same, as shown in the table.

D Effect of Refresh Threshold on Edge-CoDel

To see whether our results for Edge-CoDel were highly dependent on the refresh threshold value, consider Table 4 which shows the average FCT and RTT values for varying refresh thresholds. We find that there are very minor differences in the results as we vary this threshold, because the dominating cause for refreshing the interval is when a packet sees a queuing delay less than the CoDel target. However, the general trend is that increasing the refresh threshold increases the FCT and decreases the RTT. This is because with increasing refresh threshold, the interval is reset to the larger 100ms value less frequently. This results in more packet drops for the long flows, causing an increase in FCTs, but a decrease in the RTT values.