



# Information-Agnostic Flow Scheduling for Commodity Data Centers

Wei Bai, Li Chen, and Kai Chen, *The Hong Kong University of Science and Technology*;  
Dongsu Han, *Korea Advanced Institute of Science and Technology (KAIST)*; Chen Tian,  
*Nanjing University*; Hao Wang, *The Hong Kong University of Science and Technology*

<https://www.usenix.org/conference/nsdi15/technical-sessions/presentation/bai>

**This paper is included in the Proceedings of the  
12th USENIX Symposium on Networked Systems  
Design and Implementation (NSDI '15).**

**May 4–6, 2015 • Oakland, CA, USA**

ISBN 978-1-931971-218

**Open Access to the Proceedings of the  
12th USENIX Symposium on  
Networked Systems Design and  
Implementation (NSDI '15)  
is sponsored by USENIX**

# Information-Agnostic Flow Scheduling for Commodity Data Centers

Wei Bai<sup>1</sup>, Li Chen<sup>1</sup>, Kai Chen<sup>1</sup>, Dongsu Han<sup>2</sup>, Chen Tian<sup>3</sup>, Hao Wang<sup>1</sup>

<sup>1</sup>*SING Group @ HKUST*    <sup>2</sup>*KAIST*    <sup>3</sup>*Nanjing Univ.*

## Abstract

Many existing data center network (DCN) flow scheduling schemes minimize flow completion times (FCT) based on prior knowledge of flows and custom switch functions, making them superior in performance but hard to use in practice. By contrast, we seek to minimize FCT with no prior knowledge and existing commodity switch hardware.

To this end, we present PIAS<sup>1</sup>, a DCN flow scheduling mechanism that aims to minimize FCT by mimicking Shortest Job First (SJF) on the premise that flow size is not known *a priori*. At its heart, PIAS leverages multiple priority queues available in existing commodity switches to implement a Multiple Level Feedback Queue (MLFQ), in which a PIAS flow is gradually demoted from higher-priority queues to lower-priority queues based on the number of bytes it has sent. As a result, short flows are likely to be finished in the first few high-priority queues and thus be prioritized over long flows in general, which enables PIAS to emulate SJF without knowing flow sizes beforehand.

We have implemented a PIAS prototype and evaluated PIAS through both testbed experiments and ns-2 simulations. We show that PIAS is readily deployable with commodity switches and backward compatible with legacy TCP/IP stacks. Our evaluation results show that PIAS significantly outperforms existing information-agnostic schemes. For example, it reduces FCT by up to 50% and 40% over DCTCP [11] and L2DCT [27] respectively; and it only has a 4.9% performance gap to an ideal information-aware scheme, pFabric [13], for short flows under a production DCN workload.

## 1 Introduction

There has been a virtually unanimous consensus in the community that one of the most important goals for data center network (DCN) transport designs is to minimize the flow completion times (FCT) [11–13, 22, 26, 27].

<sup>1</sup>PIAS, Practical Information-Agnostic flow Scheduling, was first introduced in an earlier workshop paper [14] which sketched a preliminary design and the initial results.

This is because many of today’s cloud applications, such as web search, social networking, and retail recommendation, have very demanding latency requirements, and even a small delay can directly affect application performance and degrade user experience [11, 27].

To minimize FCT, most recent proposals [13, 22, 26, 34] assume prior knowledge of accurate per-flow information, e.g., flow sizes or deadlines, to achieve superior performance. For example, PDQ, pFabric and PASE [13, 22, 26] all assume flow size is known *a priori*, and attempt to approximate Shortest Job First (SJF, pre-emptive), which is the optimal scheduling discipline for minimizing the average FCT over a single link. In this paper, we question the validity of this assumption, and point out that, for many applications, such information is difficult to obtain, and may even be unavailable (§2). Existing transport layer solutions with this assumption is therefore very hard to implement in practice.

We take one step back and ask: without prior knowledge of flow size information, what is the best scheme that minimizes FCT with existing commodity switches?

Motivated by the above question, we list our key design goals as follows:

- **Information-agnostic:** Our design must not assume *a priori* knowledge of flow size information being available from the applications.
- **FCT minimization:** The solution must be able to enforce an optimal information-agnostic flow scheduling. It should minimize average and tail FCTs for latency-sensitive short flows, while not adversely affecting the FCTs of long flows.
- **Readily-deployable:** The solution must work with existing commodity switches in DCNs and be backward compatible with legacy TCP/IP stacks.

When exploring possible solution spaces, we note that some existing approaches such as DCTCP, HULL, L2DCT, etc [11, 12, 27] reduce FCT without relying on flow size information. They generally improve FCT by maintaining low queue occupation through mechanisms like adaptive congestion control, ECN, pacing, etc. However, they do not provide a full answer to our question, because they mainly perform end-host based rate control which is ineffective for flow scheduling [13, 26].

In this paper, we answer the questions with PIAS, a practical information-agnostic flow scheduling that minimizes the FCT in DCNs. In contrast to previous FCT minimization schemes [13, 22, 26] that emulate SJF by using prior knowledge of flow sizes, PIAS manages to mimic SJF with no prior information. At its heart, PIAS leverages multiple priority queues available in existing commodity switches to implement a Multiple Level Feedback Queue (MLFQ), in which a PIAS flow is gradually demoted from higher-priority queues to lower-priority queues based on the bytes it has sent during its lifetime. In this way, PIAS ensures in general that short flows are prioritized over long flows, effectively emulating SJF without knowing flow sizes beforehand.

However, we face several concrete challenges to make PIAS truly effective. First, how to determine the demoted threshold for each queue of MLFQ? Second, as flow size distribution varies across time and space, how to keep PIAS's performance in such a dynamic environment? Third, how to ensure PIAS's compatibility with legacy TCP/IP stacks in production DCNs?

For the first challenge, we address it by deriving a set of optimal demotion thresholds for MLFQ through solving a FCT minimization problem. We further show that the derived threshold setting is robust to a reasonable range of traffic distributions. This encourages us to distribute the thresholds to the end hosts for packet tagging while only performing strict priority queueing, a built-in function, in the PIAS switches.

For the second challenge, as one set of demotion thresholds works the best for a certain range of traffic distributions, PIAS adjusts thresholds to keep up with traffic dynamics. However, the key problem is that a mismatch between thresholds and underlying traffic is inevitable. Once that happens, short flows may be adversely affected by large ones in a queue, impacting on their latency. Inspired by ideas from ECN-based rate control [11], PIAS employs ECN to mitigate our mismatch problem. Our reasoning is that, by maintaining low queue occupation, short flows always see small queues and thus will not be seriously delayed even if they are mistakenly placed in a queue with a long flow due to mismatched thresholds.

For the third challenge, we employ DCTCP-like transport at the PIAS end hosts and find that PIAS interacts favorably with DCTCP or other legacy TCP protocols with ECN enabled. A potential problem is that many concurrent short flows may starve a coexisting long flow, triggering TCP timeouts and degrading application performance. We measure the extent of starvation on our testbed with a realistic workload and analyze possible solutions. We ensure that all mechanisms in PIAS can be implemented by a shim layer over NIC without touching the TCP stack.

We have implemented a PIAS prototype (§4). On the

end host, we implement PIAS as a kernel module in Linux, which resides between the Network Interface Card (NIC) driver and the TCP/IP stack as a shim layer. It does not touch any TCP/IP implementation that natively supports various OS versions. In virtualized environments, PIAS can also support virtual machines well by residing in hypervisor (or Dom 0). On the switch, PIAS only needs to enable priority queues and ECN which are both built-in functions readily supported by existing commodity switch hardware.

We evaluate PIAS on a small-scale testbed with 16 Dell servers and a commodity Pronto-3295 Gigabit Ethernet switch (Broadcom BCM#56538). In our experiments, we find that PIAS reduces the average FCT for short flows by  $\sim 37\text{-}47\%$  and  $\sim 30\text{-}45\%$  compared to DCTCP under two realistic DCN traffic patterns. It also improves the query performance by  $\sim 28\text{-}30\%$  in a Memcached [7] application (§5.1). We further dig into different design components of PIAS such as queues, optimal demotion threshold setting, ECN, and demonstrate the effectiveness of each of their contributions to PIAS's performance (§5.2).

To complement our small-scale testbed experiments, we further conduct large-scale simulations in a simulated 10/40G network using ns-2 [10]. In our simulations, we show that PIAS outperforms all existing information-agnostic solutions under realistic DCN workloads, reducing the average FCT for short flows by up to 50% and 40% compared to DCTCP and L2DCT respectively. In addition, our results show that PIAS, as a readily-deployable information-agnostic scheme, also delivers a comparable performance to a clean-slate information-aware design, pFabric [13], in certain scenarios. For example, there is only a 4.9% gap to pFabric for short flows in a data mining workload [21] (§5.3).

To make our work easy to reproduce, we made our implementation and evaluation scripts available online at: <http://sing.cse.ust.hk/projects/PIAS>.

## 2 Motivation

To motivate our design, we introduce a few cases in which the accurate flow size information is hard to obtain, or simply not available.

**HTTP chunked transfer:** Chunked transfer has been supported since HTTP 1.1 [20], where dynamically generated content is transferred during the generation process. This mode is widely used by datacenter applications. For example, applications can use chunked transfer to dump database content into OpenStack Object Storage [9]. In chunked transfer, a flow generally consists of multiple chunks, and the total flow size is not available at the start of the transmission.

**Database query response:** Query response in database systems, such as Microsoft SQL Server [8], is another example. Typically, SQL servers send partial query results as they are created, instead of buffering the result until the end of the query execution process [8]. The flow size again is not available at the start of a flow.

**Stream processing:** Stream processing systems are currently gaining popularity. In Apache Storm [2], after the master node distributes tasks to worker nodes, workers will analyze the tasks and pre-establish persistent connections with related worker nodes. During the data processing, data tuples completed in one node are continuously delivered to the next node in the stream processing chain. The amount of data to be processed is unknown until the stream finishes.

**Practical limitations:** We note that there are certain cases where the flow size information can be obtained or inferred. For example, in Hadoop [5] the mapper will first write the intermediate data to disk before the corresponding reducer starts to fetch the data, thus the flow size can be obtained in advance [28]. Even so, practical implementation issues are still prohibitive. First, we need to patch all modules in every application that generate network traffic, which is a burden for applications programmers and/or network operators. Second, current operating systems lack appropriate interface for delivering the flow size information to the transport layer. Thus, kernel modifications are also required.

### 3 The PIAS Design

#### 3.1 Design Rationale

Compared to previous solutions [13, 22, 26] that emulate SJF based on prior knowledge of flow sizes, PIAS distinguishes itself by emulating SJF with no prior information. At its core, PIAS exploits multiple priority queues available in commodity switches and implements a MLFQ, in which a PIAS flow is demoted from higher-priority queues to lower-priority queues dynamically according to its bytes sent. Through this way, PIAS enables short flows to finish in the first few priority queues, and thus in general prioritizes them over long flows, effectively mimicking SJF without knowing the flow sizes.

We note that scheduling with no prior knowledge is known as non-clairvoyant scheduling [25]. Least Attained Service (LAS) is one of the best known algorithms that minimize the average FCT in this case [30]. LAS tries to approximate SJF by guessing the remaining service time of a job based on the service it has attained so far. LAS is especially effective in DCN environments where traffic usually exhibits long-tail distribution—most flows are short and a small percent are very large [11, 21].

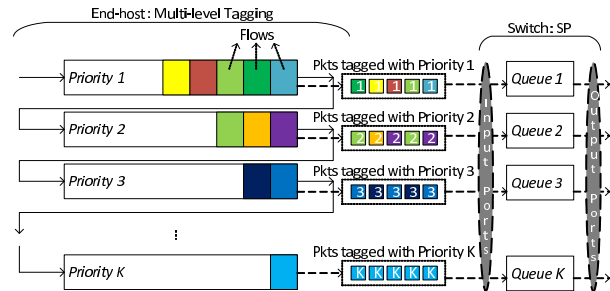


Figure 1: PIAS overview

PIAS is partially inspired by LAS. However, we find that directly enabling LAS on switches requires us to compare the amount of bytes transferred for each flow, which is not supported in existing commodity switches. Furthermore, although DCN traffic distribution is generally long-tailed, it varies across both time and space, and on some switch ports the distribution may temporarily not be so. Blindly using LAS will exacerbate the problem when multiple long flows coexist on a port, as pure LAS favors short flows but performs badly when a long flow meets a longer flow, causing the longer one to starve.

To this end, PIAS leverages multiple priority queues available in existing commodity switches (typically 4–8 queues per port [13]) to implement a MLFQ (see Figure 1). Packets in different queues of MLFQ are scheduled with strict priority, while packets in the same queue are scheduled based on FIFO. In a flow’s lifetime, it is demoted dynamically from  $i$ th queue down to the  $(i + 1)$ th queue after transmitting more bytes than queue  $i$ ’s demotion threshold, until it enters the last queue. To further prevent switches from maintaining the per-flow state, PIAS distributes packet priority tagging (indicating a flow’s sent size) to end hosts, allowing the PIAS switches to perform strict priority queuing only, which is already a built-in function in today’s commodity switches.

By implementing MLFQ, PIAS gains two benefits. First, it prioritizes short flows over large ones because short flows are more likely to finish in the first few higher priority queues while large flows are eventually demoted to lower priority queues. This effectively enables PIAS to approximate SJF scheduling that optimizes average FCT while being readily implementable with existing switch hardware. Second, it allows large flows that are demoted to the same low priority queues to share the link fairly. This helps to minimize the response time of long flows, mitigating the starvation problem.

However, there are several concrete challenges to consider in order to make PIAS truly effective. First, how to determine the demotion threshold for each queue of MLFQ to minimize the FCT? Second, as DCN traffic varies across both time and space, how to make PIAS perform efficiently and stably in such a dynamic environ-

ment? Third, how to ensure PIAS’s compatibility with legacy TCP/IP stacks in production DCNs? Next, we explain the details of the mechanism we design to address all these challenges.

### 3.2 Detailed Mechanisms

At a high level, PIAS’s main mechanisms include distributed packet tagging, switch design, and rate control.

#### 3.2.1 Packet Tagging at End-hosts

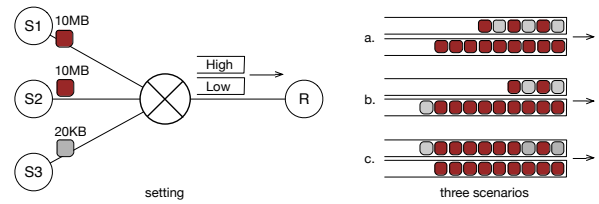
PIAS performs distributed packet tagging at end hosts as shown in Figure 1. There are  $K$  priorities  $P_i, 1 \leq i \leq K$  and  $(K - 1)$  demotion thresholds  $\alpha_j, 1 \leq j \leq K - 1$ . We assume  $P_1 > P_2 \dots > P_K$  and  $\alpha_1 \leq \alpha_2 \dots \leq \alpha_{K-1}$ .

At the end host, when a new flow is initialized, its packets will be tagged with the highest priority  $P_1$ , giving it the highest priority in the network. As more bytes are sent, the packets of this flow will be tagged with decreasing priorities  $P_j (2 \leq j \leq K)$  and enjoy decreasing priorities in the network<sup>2</sup>. The threshold to demote priority from  $P_{j-1}$  to  $P_j$  is  $\alpha_{j-1}$ .

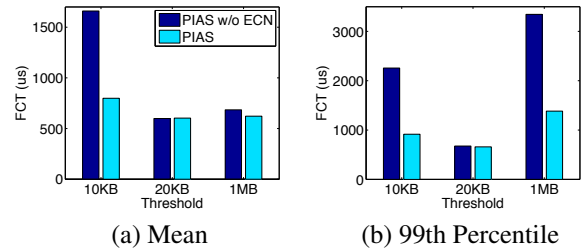
One challenge is to determine the demotion threshold for each priority to minimize the average FCT. By solving a FCT minimization problem, we derive a set of analytical solutions for optimal demotion thresholds (details in §3.3). Note that in PIAS we calculate the thresholds based on traffic information from the entire DCN and distribute the same threshold setting to all the end hosts. Our experiments and analysis show that such threshold setting is effective and also robust to a certain range of traffic variations (§5). This is a key reason we can decouple packet tagging from switches to end hosts while still maintaining good performance, which relieves the PIAS switches of having to keep the per-flow state.

As traffic changes over time, PIAS need to adjust the demotion thresholds accordingly. To keep track of traffic variations, each end host can periodically report its local traffic information to a central entity for statistics and many existing techniques exist for this purpose [37]. However, historical traffic may not reflect the future perfectly. Mismatches between threshold setting and underlying traffic are inevitable, which can hurt latency sensitive short flows. Therefore, mitigating the impact of the mismatch is a must for PIAS to operate in the highly dynamic DCNs. Our solution, as shown subsequently, is to employ ECN.

<sup>2</sup>In certain cases, applications may build persistent TCP connections to keep delivering request-response short messages for a long time. These persistent connections will eventually be assigned to the lowest priority due to the large cumulative size of bytes sent. To address this, we can periodically reset flow states based on more behaviors of traffic. For example, when a flow idles or keeps a very low average throughput for some time, we may reset the bytes sent from this flow to 0.



**Figure 2: Illustration example: (a) threshold right; (b) threshold too small, packets of short flow get delayed by long flow after prematurely demoted to the low priority queue; (c) threshold too large, packets of large flow stay too long in the high priority queue, affecting short flow.**



**Figure 3: Completion time of 20KB short flows**

#### 3.2.2 Switch Design

The PIAS switches enable the following two basic mechanisms, which are built-in functions for existing commodity switches [26].

- Priority scheduling: Packets are dequeued based on their priorities strictly when a fabric port is idle.
- ECN marking: The arriving packet is marked with Congestion Experienced (CE) if the instant buffer occupation is larger than the marking threshold.

With priority scheduling at the switches and packet tagging at the end hosts, PIAS performs MLFQ-based flow scheduling on the network fabric with stateless switches. Packets with different priority tags are classified into different priority queues. When the link is idle, the head packet from the highest non-empty priority queue is transmitted.

One may wonder why weighted fair queueing (WFQ) is not used to avoid starvation for long-lived flows. However, we choose priority scheduling for two reasons. First, priority queueing can provide better in-network prioritization and potentially achieve lower FCT than WFQ. Second, WFQ may cause packet out-of-order problem, thus degrading TCP performance.

Our intention to employ ECN is to mitigate the effect of the mismatch between the demotion thresholds and the traffic distribution. We use a simple example to illustrate the problem and the effectiveness of our solution. We connect 4 servers to a Gigabit switch as in Figure 2. One server is receiver (R) and the other three are senders (S1, S2 and S3). In our experiment, the receiver R continuously fetches 10MB data from S1 and S2, and 20KB data from S3. We configure the strict priority queueing

with 2 queues on the switch egress port to R. Since there are two priorities, we only have one demotion threshold from the high priority queue to the low priority queue. For this case, the optimal demotion threshold should be 20KB (achieving SJF in effect).

We intentionally apply three different thresholds 10KB, 20KB and 1MB, and measure the FCT of the 20KB short flows. Figure 3 shows the results of PIAS and PIAS without ECN. When the threshold is 20KB, both PIAS and PIAS without ECN achieve an ideal FCT. However, with a larger threshold (1MB) or a smaller threshold (10KB), PIAS shows obvious advantages over PIAS without ECN at both the average and 99th percentile. This is because, if the threshold is too small, packets of short flows prematurely enter the low priority queue and experience queueing delay behind long flows (see scenario (b)); if the threshold is too large, packets of long flows over-stay in the high priority queue, also affecting the latency of short flows (see scenario (c)).

By employing ECN, we can keep low buffer occupation and minimize the impact of long flows on short flows, which makes PIAS more robust to the mismatch between the demotion thresholds and traffic distribution.

### 3.2.3 Rate Control

PIAS employs DCTCP [11] as end host transport, and other legacy TCP protocols with ECN enabled can also be integrated into PIAS. We require PIAS to interact smoothly with the legacy TCP stack. One key issue is to handle flow starvation: when packets of a large flow get starved in a low priority queue for long time, this may trigger TCP timeouts and retransmissions. The frequent flow starvation may disturb the transport layer and degrade application performance. For example, a TCP connection which is starved for long time may be terminated unexpectedly.

To address the problem, we first note that PIAS can well mitigate the starvation between long flows, because two long flows in the same low priority queue will fairly share the link in a FIFO manner. In this way, PIAS minimizes the response time of each long flow, effectively eliminating TCP timeouts.

However, it is still possible that many concurrent short flows will starve a long flow, triggering its TCP timeouts. To quantify the effect, we run the web search benchmark traffic [11] (Figure 5) at 0.8 load in our 1G testbed, which has 16 servers connected to a Gigabit switch. We set RTOMin to 10ms and allocate 8 priority queues for PIAS. This experiment consists of 5,000 flows (around 5.7 million MTU-sized packets). We enable both ECN and dynamic buffer management in our switch. Hence, TCP timeouts are mainly caused by starvation rather than packet drops. We measure the number of TCP timeout-

s to quantify the extent of the starvation. We find that there are only 200 timeout events and 31 two consecutive timeout events in total. No TCP connection is terminated unexpectedly. The results indicate that, even at a high load, flow starvation is not common and will not degrade application performance adversely. We believe one possible reason is that the per-port ECN we used (see §4.1.2) may mitigate starvation by pushing back high priority flows when many packets from low priority long flows get starved. Another possible solution for handling flow starvation is treating a long-term starved flow as a new flow. For example, if a flow experiences two consecutive timeouts, we set its bytes sent back to zero. This ensures that a long flow can always make progress after timeouts. Note that the implementation of the above mechanism can be integrated to our packet tagging module without any changes to the networking stack.

Note that PIAS is free of packet reordering. This is because, during its lifetime, a PIAS flow is always demoted from a higher priority queue to a lower priority queue. In this way, an earlier packet is guaranteed to dequeue before a latter packet at each hop.

### 3.2.4 Discussion

**Local decision:** The key idea of PIAS is to emulate SJF which is optimal to minimize average FCT over a single link. However, there does not exist an optimal scheduling policy to schedule flows over an entire DCN with multiple links [13]. In this sense, similar to pFabric [13], PIAS also makes switch local decisions. This approach in theory may lead to some performance loss over the fabric [26]. For example, when a flow traverses multiple hops and gets dropped at the last hop, it causes bandwidth to be wasted on the upstream links that could otherwise have been used to schedule other flows. We note that some existing solutions [22, 26] leverage arbitration, where a common network entity allocates rates to each flow based on global network visibility, to address this problem. However, it is hard to implement because it requires non-trivial switch changes [22] or a complex control plane [26], which is against our design goal. Fortunately, local-decision based solutions maintain very good performance for most scenarios [13] and only experience performance loss at extremely high loads, e.g., over 90% [26]. However, most DCNs operate at moderate loads, e.g., 30% [16]. Our ns-2 simulation (§5.3) with production DCN traffic further confirms that PIAS works well in practice.

**Demotion threshold updating:** By employing ECN, PIAS can effectively handle the mismatch between the demotion thresholds and traffic distribution (see §5.2). This suggests that we do not need to frequently change our demotion thresholds which may be an overhead. In

this paper, we simply assume the demotion thresholds are updated periodically according to the network scale (which decides the time for information collection and distribution) and leave dynamic threshold updates as future work.

### 3.3 Optimizing Demotion Thresholds

In this section, we describe our formulation to derive the optimal demotion thresholds for minimizing the average FCT. We identify this problem as a Sum-of-Linear-Ratios problem and provide a method to derive the optimal thresholds analytically for *any* given load and flow size distribution. We find that the demotion thresholds depend on both load and flow size distribution. As flow size distribution and load change across both time and space, ideally, one should use different thresholds for different links at different times. However, in practice, it is quite challenging to obtain such fine-grained link level traffic information across the entire DCN. Hence, we use the overall flow size distribution and load measured in the entire DCN as an estimate to derive a common set of demotion thresholds for all end hosts. We note that this approach may not be theoretically optimal and there is room for improvement. However, it is more practical and provides considerable gains, as shown in the evaluation (§5).

**Problem formulation:** We assume there are  $K$  priority queues  $P_i (1 \leq i \leq K)$  where  $P_1$  has the highest priority. We denote the threshold for demoting the priority from  $j-1$  to  $j$  as  $\alpha_{j-1} (2 \leq j \leq K)$ . We define  $\alpha_K = \infty$ , so that the largest flows are all in this queue, and  $\alpha_0 = 0$ . The flows, indexed from  $i=1$  to  $N$ , have sizes  $x_i$ . Denote the cumulative density function of flow size distribution as  $F(x)$ , thus  $F(x)$  is the probability that a flow size is no larger than  $x$ . Note that we do not assume the any specific properties of  $F(x)$ , and it is used for convenience for the following derivation.

Let  $\theta_j = F(\alpha_j) - F(\alpha_{j-1})$ , the percentage of flows with sizes in  $[\alpha_{j-1}, \alpha_j)$ . For a flow with size in  $[\alpha_{j-1}, \alpha_j)$ , it experiences the delays in different priorities up to the  $j$ -th priority. Denote  $T_j$  as the average time spent in the  $j$ -th queue. For a flow with size  $x$ , let  $x^+$  be the residual size of this flow tagged with its lowest priority. Thus,  $x^+ = x - \alpha_{j_{max}}(x)$ , where  $\alpha_{j_{max}}(x)$  is the largest demotion threshold less than  $x$ , and let  $j_{max}(x)$  be the index of this threshold.

So the average FCT for this flow is:  $T(x) = \sum_{l=1}^{j_{max}(x)} T_l + \frac{x^+}{1 - \rho_{j_{max}(x)}}$ .

The second term is bounded by  $T_{j_{max}(x)}$ , thus an upper bound is therefore:  $T(x) \leq \sum_{l=1}^{\min(j_{max}(x)+1, K)} T_l$ .

We have the following optimization problem, where we choose an optimal set of thresholds  $\{\alpha_j\}$  to minimize

the objective: the average FCT of flows on this bottleneck link:

$$\begin{aligned} \min_{\{\alpha_j\}} \quad & \mathcal{T} = \sum_{l=1}^K (\theta_l \sum_{m=1}^l T_m) = \sum_{l=1}^K (T_l \sum_{m=l}^K \theta_m) \\ \text{subject to} \quad & \alpha_0 = 0, \alpha_K = \infty \\ & \alpha_{j-1} < \alpha_j, j=1, \dots, K \end{aligned} \quad (1)$$

**Analysis:** For convenience, we use  $\theta$ s to equivalently replace  $\alpha$ s. With a traffic load of  $\rho$ , the average time in the  $l$ th queue,  $T_l$ , can be expressed as:  $T_l = \frac{\theta_l \rho}{1 - \rho F(\alpha_{l-1})}$  (assuming M/M/1 queues<sup>3</sup>). Since  $\sum_{m=1}^l \theta_m = \sum_{m=l}^K F(\alpha_m) - F(\alpha_{m-1})$ , we can re-express the objective as:  $\mathcal{T} = \sum_{l=1}^K T_l (1 - F(\alpha_{l-1}))$ .  $\sum_{l=1}^K \theta_m = 1$ .

Since  $\sum_{l=1}^K \theta_m = 1$ , we can finally transform the problem as:

$$\max_{\{\theta_l\}} \quad \mathcal{T}'' = \sum_{l=1}^{K-1} \frac{\theta_l}{1 - \rho (\sum_{m=1}^{l-1} \theta_m)} + \frac{1 - \sum_{m=1}^{K-1} \theta_m}{1 - \rho \sum_{m=1}^{K-1} \theta_m} \quad (2)$$

which is a Sum-of-Linear-Ratios (SoLR) problem [32], a well-known class of fractional programming problems. The only constraint is that  $\theta_l \geq 0, \forall l$ . This formulation is interesting because the upperbound of average FCT is independent of the flow distribution  $F(x)$ , and only concerns the  $\theta$ s, which represent the percentages of traffic in different queues. Thus  $F(x)$  is needed only when we calculate the thresholds, so we can first obtain the optimal set of  $\theta$ s for all links, and then derive the priority thresholds based on  $F(x)$ .

**Solution method:** Generally, the SoLR problem is  $\mathcal{NP}$ -hard [17], and the difficulty in solving this class of problems lies in the lack of useful properties to exploit. For our problem, however, we find a set of properties that can lead to a closed-form analytical solution to Problem 2. We describe the derivation procedure as follows:

Consider the terms in the objective. Since the traffic load  $\rho \leq 1$ , we have  $\rho (\sum_{m=1}^{l-1} \theta_m) \leq \sum_{m=1}^{l-1} \theta_m$ . Also,  $\theta_l + \sum_{m=1}^{l-1} \theta_m = \sum_{m=1}^l \theta_m \leq 1$ . Thus we have:

$$\theta_l \leq \sum_{m=l}^K \theta_m = 1 - \sum_{m=1}^{l-1} \theta_m \leq 1 - \rho (\sum_{m=1}^{l-1} \theta_m) \quad (3)$$

The property to exploit is as follows: each term in the summation,  $\theta_l / (1 - \rho \sum_{m=1}^{l-1} \theta_m)$ , is no larger than 1, and to maximize the summation, we should make the numerator and denominator as close as possible, so that the ratio is close to 1.

Consider the first two portions,  $\theta_1$  and  $\theta_2$ . By making the numerator and denominator equal, we have:

$$\theta_2 = 1 - \rho \theta_1 \quad (4)$$

We can obtain the expression of the third portion and

<sup>3</sup> We use M/M/1 queues to simplify the analysis and obtain a closed-form solution. Similar derivation can also be conducted using M/G/1 queues [13], but the solution is very complicated.

all the portions after it iteratively:

$$\theta_3 = 1 - \rho(\theta_2 + \theta_1) = 1 - \rho(1 - (1 - \rho)\theta_1) \quad (5)$$

In this way, by the formula  $\theta_l = 1 - \rho(\sum_{m=1}^{l-1} \theta_m)$ , each portion  $\theta_l$  can be represented by an expression of  $\theta_1$  iteratively. In addition, by the constraint  $\sum_{l=1}^K \theta_l = 1$ , we can obtain the analytical expressions of all  $\theta_s$ , which represents the percentages of the traffic in different priority queues on the link. Given traffic load  $\theta_s$  and flow size distribution  $F(\cdot)$ , we can express the thresholds in the unit of bytes, and this representation is implemented at the end host.

## 4 Implementation and Testbed Setup

### 4.1 PIAS Implementation

We have implemented a prototype of PIAS. We now describe each component of our prototype in detail.

#### 4.1.1 Packet Tagging

The packet tagging module is responsible for maintaining per-flow state and marking packets with priority at the end hosts. We implement it as a kernel module in Linux. The packet tagging module resides between the TCP/IP stack and Linux TC, which consists of three components: a NETFILTER [6] hook, a hash based flow table, and a packet modifier.

The operations are as follows: 1) the NETFILTER hook intercepts all outgoing packets using the LOCAL\_OUT hook and directs them to the flow table. 2) Each flow in the flow table is identified by the 5-tuple: src/dst IPs, src/dst ports and protocol. When a packet comes in, we identify the flow it belongs to (or create a new entry) and increment the amount of bytes sent. 3) Based on the flow information, the packet modifier sets the packet's priority by modifying the DSCP field in the IP header to the corresponding value.

Offloading techniques like large segmentation offloading (LSO) may degrade the accuracy of packet tagging. With LSO, the packet tagging module may not be able to set the right DSCP value for each individual MTU-sized packet within a large segment. To quantify this, we sample more than 230,000 TCP segments with payload data in our 1G testbed and find that the average segment size is only 7,220 Bytes. This has little impact on packet tagging. We attribute this to the small window size in DCN environment which has small bandwidth-delay product and large number of concurrent connections. We expect that the final implementation solution for packet tagging should be in NIC hardware to permanently avoid this interference.

To quantify system overhead introduced by the PIAS packet tagging module, we installed it on a Dell PowerEdge R320 server with an Intel 82599EB 10GbE NIC

and measured CPU usage. LSO is enabled in this experiment. We started 8 TCP long flows and achieved  $\sim 9.4$ Gbps goodput. The extra CPU usage introduced by PIAS is  $< 1\%$  compared with the case where the PIAS packet tagging module is not enabled.

#### 4.1.2 Switch Configuration

We enforce strict priority queues at the switches and classify packets based on the DSCP field. Similar to previous work [11, 35], we use ECN marking based on the instant queue lengths with a single marking threshold. In addition to the switch queueing delay in the network, sender NIC also introduces latency because it is actually the first contention point of the fabric [13, 24]. Hardware and software queues at the end hosts can introduce large queueing delay, which might severely degrade the application performance [23, 36]. To solve this problem, our software solution hooks into the TX datapath at POST\_ROUTING and rate-limits outgoing traffic at the line rate. Then, we perform ECN marking and priority queueing at the end host as well as the switches.

**Per-queue vs per-port ECN marking:** We observe that some of today's commodity switching chips offer multiple ways to configure ECN marking when configured to use multiple queues per port. For example, our Broadcom BCM#56538-based switch allows us to enable either per-queue ECN marking or per-port ECN marking. In per-queue ECN marking, each queue has its own marking threshold and performs ECN marking independently to other queues. In per-port ECN marking, each port is assigned a single marking threshold and marks packets when the sum of all queue sizes belonging to the port exceeds the marking threshold.

Per-queue ECN is widely used in many DCN transport protocols [11, 26, 33], however, we find it has limitations when supporting multiple queues. Each queue requires a moderate ECN marking threshold  $h$  to fully utilize the link independently (e.g.,  $h=20$  packets for 1G and 65 packets for 10G in DCTCP [11]). Thus, supporting multiple queues may require the shared memory be at least multiple times (e.g., 8) the marking threshold, which is not affordable for most shallow buffered commodity switches. For example, our Pronto-3295 switch has 4MB ( $\approx 2667$  packets) memory shared by 384 queues (48x 1G ports with 8 queues per port). If we set  $h=20$  packets as suggested above, we need over 11MB memory in the worst case, otherwise when the traffic is bursty, the shallow buffer may overflow before ECN takes effect.

Per-port ECN, to the best of our knowledge, has rarely been exploited in recent DCN transport designs. Although per-port ECN marking cannot provide ideal isolation among queues as per-queue ECN marking (Figure 12), it can provide much better burst tolerance and support a larger number of queues in shallow buffered



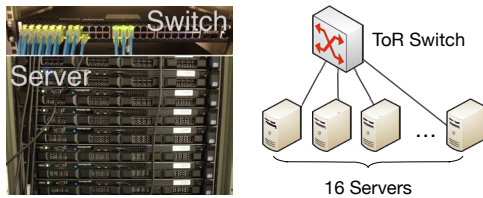


Figure 4: Testbed Topology

switches. Moreover, per-port ECN marking can potentially mitigate the starvation problem. It can push back high priority flows when many packets of low priority flows get queued in the switch. Therefore, we use per-port ECN marking.

#### 4.1.3 Rate Control

We use the open source DCTCP patch [4] for Linux 2.6.38.3. We further observe an undesirable interaction between the open-source DCTCP implementation and our switch. The DCTCP implementation does not set the ECN-capable (ECT) codepoint on TCP SYN packets and retransmitted packets, following the ECN standard [31]. However, our switch drops any non-ECT packets from ECN-enabled queues, when the instant queue length is larger than the ECN marking threshold. This problem severely degrades the TCP performance [35]. To address this problem, we set ECT on every TCP packet at the packet modifier.

## 4.2 Testbed Setup

We built a small testbed that consists of 16 servers connected to a Pronto 3295 48-port Gigabit Ethernet switch with 4MB shared memory, as shown in Figure 4. Our switch supports ECN and strict priority queuing with at most 8 class of service queues [1]. Each server is a Dell PowerEdge R320 with a 4-core Intel E5-1410 2.8GHz CPU, 8G memory, a 500GB hard disk, and a Broadcom BCM5719 NetXtreme Gigabit Ethernet NIC. Each server runs Debian 6.0-64bit with Linux 2.6.38.3 kernel. By default, advanced NIC offload mechanisms are enabled to reduce the CPU overhead. The base round-trip time (RTT) of our testbed is around 100 $\mu$ s.

In addition, we have also built a smaller 10G testbed for measuring the end host queuing delay in the high speed network. We connect three servers to the same switch (Pronto 3295 has four 10GbE ports). Each server is equipped with an Intel 82599EB 10GbE NIC.

## 5 Evaluation

We evaluate PIAS using a combination of testbed experiments and large-scale ns-2 simulations. Our evaluation centers around four key questions:

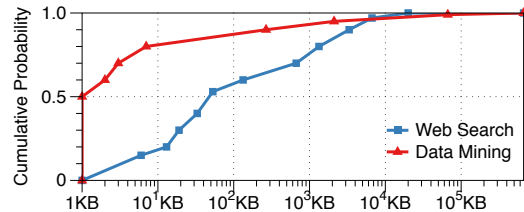


Figure 5: Traffic distributions used for evaluation.

- How does PIAS perform in practice?** Using realistic workloads in our testbed experiments, we show that PIAS reduces the average FCT of short flows by  $\sim$ 37-47% with the web search workload [11] and  $\sim$ 30-45% with the data mining workload [21] compared to DCTCP. In an application benchmark with Memcached [7], we show that PIAS achieves  $\sim$ 28-30% lower average query completion time than DCTCP.
- How effective are individual design components of PIAS, and how sensitive is PIAS to parameter settings?** We show that PIAS achieves reasonable performance even with two queues. We also demonstrate that ECN is effective in mitigating the harmful effect of a mismatch between the demotion thresholds and traffic, but PIAS performs the best with the optimal threshold setting.
- Does PIAS work well even in large datacenters?** Using large-scale ns-2 simulations, we show that PIAS scales to multi-hop topologies and performs best among all information-agnostic schemes (DCTCP [11], L2DCT [27], and LAS [30]). PIAS shows a 4.9% performance (measured in average FCT) gap from pFabric [13], an idealized information-aware scheme, for short flows in the data mining workload.
- How robust is the PIAS threshold setting?** With the same set of thresholds, we experiment with both extremely-biased and realistic traffic patterns to show the robustness of the PIAS threshold setting. We also demonstrate the optimality region and close-to-optimality region is large for PIAS analytically. (For space limitations, please see the details in our technical report [15]).

### 5.1 Testbed Experiments

**Setting:** PIAS uses 8 priority queues by default and enable per-port ECN marking as discussed in Section 4. We set the ECN marking threshold to be 30KB as DCTCP [11] recommends. Our MLFQ demotion thresholds are derived as described in Section 3.

We use two realistic workloads, a web search workload [11] and a data mining workload [21] from production datacenters. Their overall flow size distributions are shown in Figure 5. We also evaluate PIAS using an application benchmark with Memcached [7].

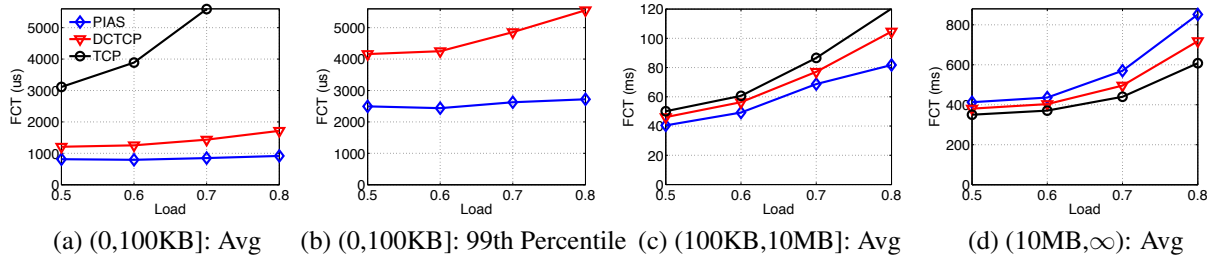


Figure 6: Web search workload: FCT across different flow sizes. TCP's performance is outside the plotted range of (b)

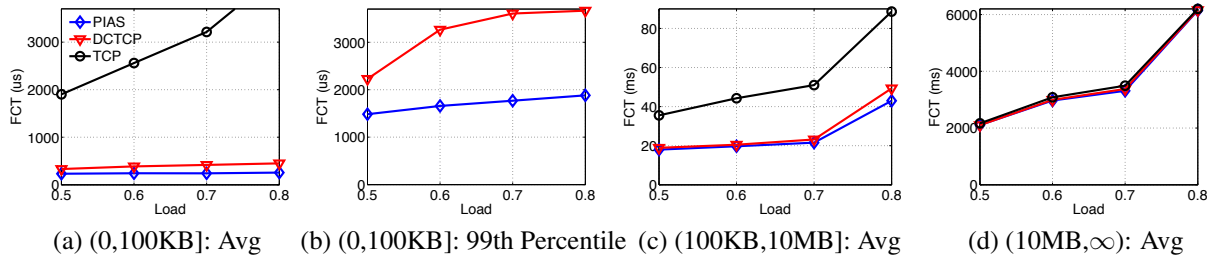


Figure 7: Data mining workload: FCT across different flow sizes. TCP's performance is outside the plotted range of (b)

**Results with realistic workloads:** For this experiment, we developed a client/server model to generate dynamic traffic according to realistic workloads and measure the FCT on application layer. The client application, running on 1 machine, periodically generates requests to the other machines to fetch data. The server applications, running on 15 other machines, respond with requested data. The requests are generated based on a Poisson process. We evaluate the performance of PIAS, DCTCP and TCP, while varying the network loads from 0.5 to 0.8. Given the average traffic load in DCNs is moderate (for example, 30% [16]), a long-term load of over 80% is less likely in practice.

Figure 6 and Figure 7 show the average FCT across small (0,100KB] (a, b), medium (100KB,10MB] (c), and large (10MB,∞) (d) flows, respectively; for the web search and data mining workloads, respectively.

We make the following three observations: First, for both workloads, PIAS achieves the best performance in both the average and 99th percentile FCTs of small flows. Compared to DCTCP, PIAS reduces the average FCT of small flows by ~37-47% for the web search workload and ~30-45% for the data mining workload. The improvement of PIAS over DCTCP in the 99th percentile FCT of short flows is even larger: ~40-51% for the web search workload and ~33-48% for the data mining workload. Second, PIAS also provides the best performance in medium flows. It achieves up to 22% lower average FCT of medium flows than DCTCP in the web search workload. Third, PIAS does not severely penalize the large flows. For example, from Figure 6 (d) and Figure 7 (d), we can see that for the data mining workload PIAS is comparable or slightly better than TCP and DCTCP, while for the web search workload it is worse than

DCTCP by over 10%. This is expected because PIAS prioritizes short flows over long flows and ~60% of all bytes in the web search workload are from flows smaller than 10MB. Note that this performance gap would not affect the overall average FCT since datacenter workloads are dominated by small and medium flows (e.g., only ~3% flows are larger than 10MB in the web search workload). In fact, PIAS achieves ~9% and ~17% lower overall average FCT than DCTCP and TCP at 0.8 load in the web search workload.

**Results with the Memcached application:** To assess how PIAS improves the performance of latency-sensitive applications, we build a Memcached [7] cluster with 16 machines. One machine is used as a client and the other 15 are used as servers to emulate a partition/aggregate soft-real-time service [11, 34]. We pre-populate server instances with 4B-key, 1KB-value pairs. The client sends a GET query to all 15 servers and each server responds with a 1KB value. A query is completed only when the client receives all the responses from the servers. We measure the query completion time as the application performance metric. Since a 1KB response can be transmitted within one RTT, the query completion time is mainly determined by the tail queueing delay. The base query completion time is around 650us in our testbed. We also generate background traffic, a mix of mice flows and elephant flows following the distribution of the web search workload [11]. We use queries per second, or *qps*, to denote the application load. We vary the load of the background traffic from 0.5 to 0.8 and compare the performance of PIAS with that of DCTCP.

Figure 8 and Figure 9 show the results of the query completion time at 20 and 40 *qps* loads respectively. S-

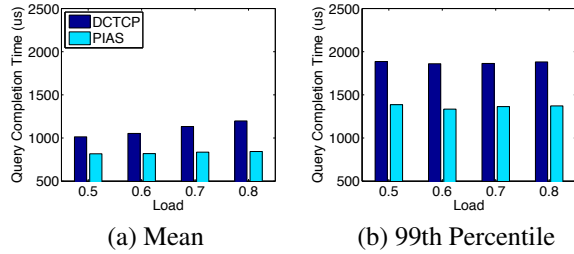


Figure 8: Query completion time at 20 qps

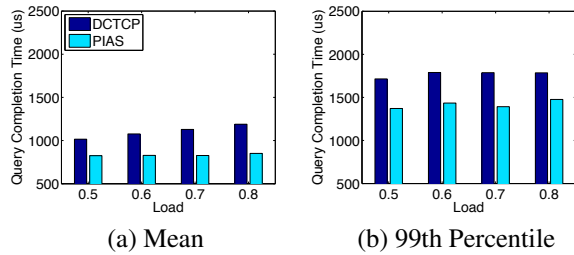


Figure 9: Query completion time at 40 qps

since we enable both dynamic buffer management and ECN on the switch, none of queries suffers from TCP timeout. With the increase in background traffic load, the average query completion time of DCTCP also increases (1016–1189us at 40qps and 1014–1198us at 20qps). By contrast, PIAS maintains a relatively stable performance. At 0.8 load, PIAS can achieve ~28-30% lower average query completion times than those of DCTCP. Moreover, PIAS also reduces the 99th percentile query completion time by ~20-27%. In summary, PIAS can effectively improve the performance of the Memcached application by reducing the queueing delay of short flows.

**End host queueing delay:** The above experiments mainly focus on network switching nodes. PIAS extends its switch design to the end hosts as the sender’s NIC is actually the first contention point of the fabric [13, 24].

To quantify this, we conduct an experiment in our 10G setting with three servers (one sender and two receivers). We start several (1 to 8) long-lived TCP flows from the sender to a receiver. Then we measure RTT from the sender to the other receiver by sending ICMP ping packets. Without PIAS, ping packets could experience up to 6748us queueing delay with 8 background flows. Then we deploy a 2-queue PIAS end host scheduling module (as described in §4.1.2) with a threshold of 100KB. Each ICMP packet is identified as a new flow by PIAS. We measure the RTTs with PIAS and compare them with the results without PIAS in Figure 10. In general, PIAS can significantly reduce the average RTT to ~200us and ensure that the 99th percentile RTT is smaller than 450us. Note that the PIAS scheduling module does not affect network utilization and large flows still maintain more than 9Gbps goodput during the experiment. Since we enable LSO to reduce CPU overhead,

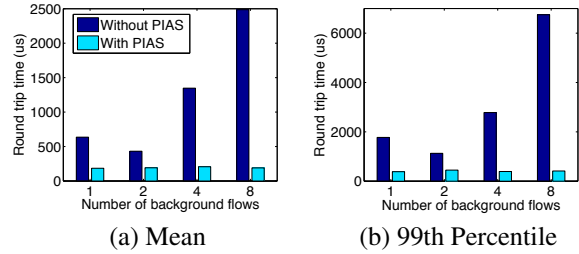


Figure 10: RTT with background flows

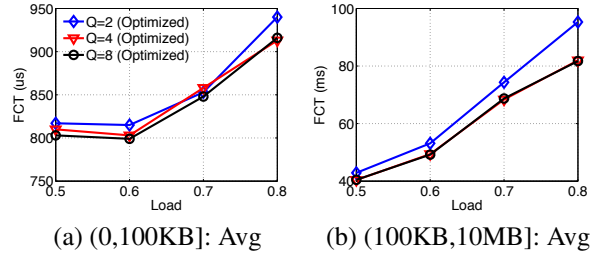


Figure 11: Web search workload with different queues

it is difficult for us to achieve fine-grained transmission control and some delay may still exist in NIC’s transmission queues. We believe there is still room to improve by offloading the scheduling to NIC hardware [29].

## 5.2 PIAS Deep Dive

In this section, we conduct a series of targeted experiments to answer the following three questions:

- **How sensitive is PIAS to the number of queues available?** Network operators may reserve some queues for other usage while some commodity switches [3] only support 2 priority queues. We find that, even with only 2 queues, PIAS still effectively reduces the FCT of short flows. However, in general, more queues can further improve PIAS’s overall performance.
- **How effective is ECN in mitigating the mismatch?** ECN is integrated into PIAS to mitigate the mismatch between the demotion thresholds and traffic distribution. In an extreme mismatch scenario, we find that without ECN, PIAS’s performance suffers with medium flows and is worse than DCTCP. However, with ECN, PIAS effectively mitigates this problem, and is better than, or at least comparable to DCTCP.
- **What is the effect of the optimal demotion thresholds?** Compared to PIAS with thresholds derived from simple heuristics, PIAS with the optimal demotion thresholds achieves up to ~10% improvement in medium flows, which improves the overall performance.

**Impact of number of queues:** In general, the more queues we use, the better we can segregate different flows, thus improving overall performance. For this

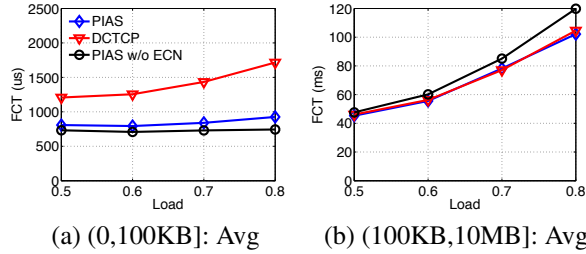


Figure 12: Web search workload with mismatch thresholds derived from data mining workload

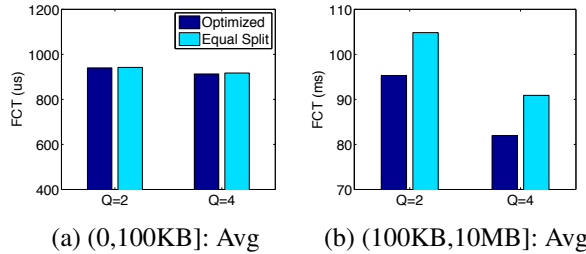


Figure 13: Web search workload with different thresholds

experiment, we generate traffic using the web search workload and do the evaluation with 2, 4 and 8 priority queues. The results are shown in Figure 11. We observe that three schemes achieve the similar average FCT of short flows. Even at 0.8 load, the FCT of 2 queues is within  $\sim 2.5\%$  of that of 8 queues. As expected, the average FCT of medium flows improves with the increasing number of queues. For example, PIAS with 4 queues and 8 queues provide similar performance, but improve the FCT by 14.3% compared to 2 queues. The takeaway is that PIAS can effectively reduce FCT of short flows even with 2 queues and more queues can further improve PIAS's overall performance.

**Effect of ECN under thresholds-traffic mismatch:** We evaluate the performance of the web search workload while using the optimal demotion thresholds derived from the data mining workload. We compare PIAS, PIAS without ECN, and DCTCP. Figure 12 shows the results of the average FCT of short and medium flows. Both PIAS and PIAS without ECN greatly outperforms DCTCP in short flows. PIAS without ECN is even slightly better than PIAS. That is because, when using per-port ECN, packets in a high priority queue may get marked due to buffer occupation in a low priority queue. However, PIAS without ECN shows the worst performance in medium flows while being obviously worse than DCTCP. This is because, due to the mismatch between demoting thresholds and traffic distribution, medium flows and large flows coexist in the lower priority queues. Without ECN, packets from medium flows would experience queuing delay behind large flows. With ECN, PIAS effectively mitigates this side-effect by keeping low buffer occupation as explained in §3.2.2.

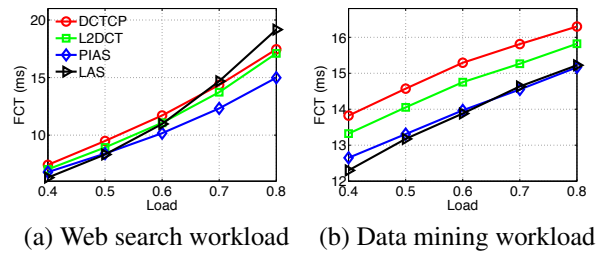


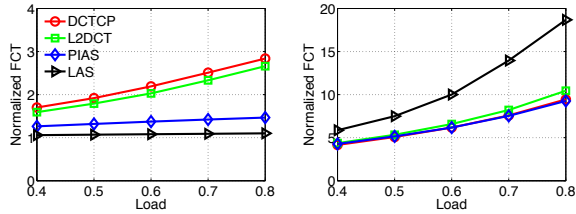
Figure 14: Overall average flow completion time

**Impact of demotion thresholds:** To explore the effectiveness of our optimal demotion threshold setting, we compare the optimized PIAS with the PIAS using thresholds derived from the equal split heuristic as [13]. More specifically, given a flow size distribution and  $N$  queues, we set the first threshold to the size of  $100/N$ th percentile flow, the second threshold to the size of  $200/N$ th percentile flow, and so on. We run the web search workload at 0.8 load and summarize results in the Figure 13. We test PIAS with 2 and 4 queues. We observe that there is an obvious improvement in the average FCT of medium flows with the optimized thresholds. Specifically, PIAS (4-queue) with the optimized thresholds can achieve  $\sim 10\%$  lower FCT for medium flows than that of equal split, and a 9% improvement for the 2-queue PIAS. This partially validates the effectiveness of our optimal threshold setting. We further conduct deep analysis on this in [15].

### 5.3 Large-scale NS-2 Simulations

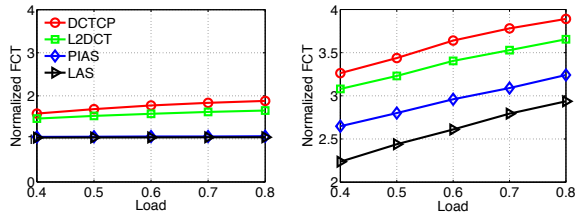
We use ns-2 [10] simulations to answer three questions.

- **How does PIAS perform compared to information-agnostic schemes?** PIAS outperforms DCTCP [11] and L2DCT [27] in general, and significantly improves their average FCTs for short flows by 50% and 40% respectively. Furthermore, PIAS is close to LAS for short flows and greatly outperforms LAS for long flows, reducing its average FCT by 50% in the web search workload.
- **How does PIAS perform compared to information-aware schemes?** As a practical information-agnostic scheme, PIAS can also deliver comparable performance to a clean-slate information-aware design, p-Fabric [13], in certain scenarios. For example, it only has a 4.9% gap to pFabric for short flows in the data mining workload.
- **How does PIAS perform in the oversubscribed network?** In a 3:1 oversubscribed topology with ECMP load balancing, PIAS still delivers very good performance. Compared to DCTCP, the average FCT for the short flows with PIAS is  $\sim 26\%$  lower in the web search workload.



(a) (0,100KB]: Avg (b) (10MB,∞): Avg

Figure 15: Web search workload: Normalized FCT



(a) (0,100KB]: Avg (b) (10MB,∞): Avg

Figure 16: Data mining workload: Normalized FCT

**Setting:** We use a leaf-spine topology with 9 leaf (ToR) switches to 4 spine (Core) switches. Each leaf switch has 16 10Gbps downlinks (144 hosts) and 4 40Gbps uplinks to the spine, forming a non-oversubscribed network. The base end-to-end round-trip time across the spine (4 hops) is  $85.2\mu\text{s}$ . We use packet spraying [19] for load balancing and disable dupACKs to avoid packet reordering. Again, we use the web search and data mining workloads as above.

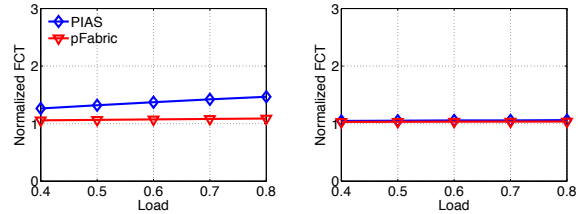
### 5.3.1 Comparison with Information-agnostic Schemes

We mainly compare PIAS with three other information-agnostic schemes: DCTCP, L2DCT [27] and LAS [30].

**Overall performance:** Figure 14 shows the average FCT of information-agnostic schemes under different workloads and load levels. From the figure, we see that PIAS performs well overall. First, PIAS has an obvious advantage over DCTCP and L2DCT in all cases. Second, PIAS is close to LAS in the data mining workload, and significantly outperforms LAS by 28% (at 0.8 load) in the web search workload. This is because PIAS effectively mitigates the starvation between long flows unlike LAS. In the data mining workload, there are not so many large flows on the same link concurrently. As a result, LAS does not suffer from starvation as significantly.

**Breakdown by flow size:** We now breakdown the average FCT across different flow sizes, (0, 100KB] and (10MB, ∞) (Figure 15 and 16). We normalize each flow’s actual FCT to the best possible value it can achieve in the idle fabric.

For short flows in (0,100KB], we find that PIAS significantly outperforms both DCTCP and L2DCT, improv-



(a) Web search workload (b) Data mining workload

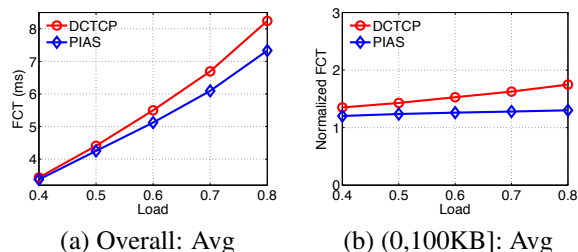
Figure 17: Average normalized FCT for (0,100KB]

ing the average FCT by up to 50% and 40% respectively. This is because DCTCP and L2DCT use reactive rate control at the end hosts, which is not as effective as PIAS for in-network flow scheduling. We further observe that PIAS achieves similar performance as LAS for short flows. PIAS only performs slightly worse than LAS in the web search workload when there is a packet drop or an explicit congestion notification.

For long flows in (10MB,∞), we find that PIAS is slightly worse than LAS in data mining workload, but performs significantly better in the web search workload (50% reduction in FCT at 0.8 load). This is because, in the web search workload, it is common that multiple large flows are present in the same link. In such scenarios, LAS always stops older flows to send new flows. Since large flows usually take a very long time to complete, it causes a serious starvation problem. However, with PIAS, large flows receive their fair sharing in the lowest priority queue, which mitigates this problem. Furthermore, PIAS performs similarly to DCTCP under the web search workload and achieves 17% lower FCT in the data mining workload.

### 5.3.2 Comparison with Ideal Information-aware Schemes

We compare PIAS to an ideal information-aware approach for DCN transport, pFabric [13], on small flows of the two workloads. We note that the most recent work PASE [26] can achieve better performance than pFabric in particular scenarios (e.g., very high load and single rack). However in our topology setting with realistic workloads, pFabric is better than PASE and PDQ [22], and achieves near-optimal performance. Thus, we directly compare PIAS with pFabric. The result is shown in Figure 17. In general, PIAS delivers comparable average FCT for short flows as pFabric, particularly within 4.9% in the data mining workload. We find that the gap between PIAS and pFabric is smaller in the data mining workload than that in the web search workload. This is mainly due to the fact that the data mining workload is more skewed than the web search workload. Around 82% flows in the data mining are smaller than 100KB, while only 54% of flows in the web search are small-



**Figure 18: Web search workload on a 3:1 oversubscribed topology with ECMP load balancing.**

er than 100KB. For the web search workload, it is more likely that large flows coexist with short flows in the high priority queues temporarily, increasing the queueing delay for short flows. pFabric, by assuming prior knowledge of flow sizes, is immune to such problem.

### 5.3.3 Performance in oversubscribed network

Finally, we evaluate PIAS on a 3:1 oversubscribed network with ECMP load balancing. In this topology, there are 3 leaf switches and 4 spine switches. Each leaf switch is connected to 48 hosts with 10Gbps links and 4 spine switches with 40Gbps links. Given that the source and destination of each flow is generated randomly, one-third of traffic is intra-ToR and the rest is inter-ToR traffic. Hence, the load at the fabric’s core is twice the load at the edge. We repeat the web search workload and compare PIAS with DCTCP. Figure 18 gives the results. Note that the load in the figure is at the fabric’s core. Compared to DCTCP, PIAS achieves up to ~26% and ~11% lower average FCT for short flows and all the flows respectively.

## 6 Related Work

We classify previous work on minimizing FCT in DCNs into two categories: information-agnostic solutions (e.g., [11, 12, 27]) and information-aware solutions (e.g., [13, 22, 26]).

Information-agnostic solutions [11, 12, 27] generally improve the FCT for short flows by keeping low queue occupancy. For example, DCTCP [11] tries to keep the fabric queues small by employing an ECN-based adaptive congestion control algorithm to throttle long elephant flows. L2DCT [27] adds bytes sent information to DCTCP [11]. HULL [12] further improves the latency of DCTCP by trading network bandwidth. In summary, these solutions mainly perform end-host based rate control which is ineffective for flow scheduling. By contrast, PIAS leverages in-network priority queues to emulate SJF for flow scheduling, which is more efficient in terms of FCT minimization.

Information-aware solutions [13, 22, 26] attempt to approximate ideal Shortest Remaining Processing Time (SRPT) scheduling. For example, PDQ [22] employs switch arbitration and uses explicit rate control for flow scheduling. pFabric [13] decouples flow scheduling from rate control and achieves near-optimal FCT with decentralized in-network prioritization. PASE [26] synthesizes the strengths of previous solutions to provide good performance. In general, these solutions can potentially provide ideal performance, but they require non-trivial modifications on switch hardware or a complex control plane for arbitration. By contrast, PIAS does not touch the switch hardware or require any arbitration in the control plan, while still minimizing FCT.

There are also some other efforts [18, 33, 34] targeting at meeting flow deadlines. D3 [34] assigns rates to flows according to their sizes and deadlines explicitly, whereas D2TCP [33] and MCP [18] add deadline-awareness to ECN-based congestion window adjustment implicitly. They all require prior knowledge of flow information and do not directly minimize FCT, unlike PIAS.

## 7 Conclusion

Through PIAS, we leverage existing commodity switches in DCNs to minimize the average FCT for flows, especially the smaller ones, without assuming any prior knowledge of flow sizes. We have implemented a PIAS prototype using all commodity hardware and evaluated PIAS through a series of small-scale testbed experiments as well as large-scale packet-level ns-2 simulations. Both our implementation and evaluation results show that PIAS is a viable solution that achieves all our design goals.

## Acknowledgements

This work was supported by the Hong Kong RGC ECS 26200014, the China 973 Program under Grant No. 2014CB340303, the ICT R&D program of MSIP/IITP, Republic of Korea [14-911-05-001], and the Basic Science Research Program of NRF funded by MSIP, Rep. of Korea (2013R1A1A1076024). We thank Mohammad Alizadeh for sharing codes of pFabric, Haitao Wu for insightful discussions, our shepherd Dushyanth Narayanan and the anonymous NSDI reviewers for their constructive comments.

## References

- [1] <http://www.pica8.com/documents/pica8-datasheet-picos.pdf>.
- [2] “Apache Storm,” <https://storm.incubator.apache.org/>.

- [3] “Cisco Nexus 5500 Series NX-OS Quality of Service Configuration Guide,” [http://www.cisco.com/c/en/us/td/docs/switches/datacenter/nexus5500/sw/qos/7x/b\\_5500\\_QoS.Config\\_7x.pdf](http://www.cisco.com/c/en/us/td/docs/switches/datacenter/nexus5500/sw/qos/7x/b_5500_QoS.Config_7x.pdf).
- [4] “DCTCP Patch,” <http://simula.stanford.edu/~alizade/Site/DCTCP.html>.
- [5] “Hadoop,” <http://hadoop.apache.org/>.
- [6] “Linux netfilter,” <http://www.netfilter.org/>.
- [7] “Memcached,” <http://memcached.org/>.
- [8] “Microsoft SQL Server,” <http://www.microsoft.com/en-us/server-cloud/products/sql-server/>.
- [9] “OpenStack Object Storage,” <http://docs.openstack.org/api/openstack-object-storage/1.0/content/chunked-transfer-encoding.html>.
- [10] “The Network Simulator NS-2,” <http://www.isi.edu/nsnam/ns/>.
- [11] M. Alizadeh, A. Greenberg, D. A. Maltz, J. Padhye, P. Patel, B. Prabhakar, S. Sengupta, and M. Sridharan, “Data center TCP (DCTCP),” in *SIGCOMM 2010*.
- [12] M. Alizadeh, A. Kabbani, T. Edsall, B. Prabhakar, A. Vahdat, and M. Yasuda, “Less is more: trading a little bandwidth for ultra-low latency in the data center,” in *NSDI 2012*.
- [13] M. Alizadeh, S. Yang, M. Sharif, S. Katti, N. McKeown, B. Prabhakar, and S. Shenker, “pFabric: Minimal near-optimal datacenter transport,” in *SIGCOMM 2013*.
- [14] W. Bai, L. Chen, K. Chen, D. Han, C. Tian, and W. Sun, “PIAS: Practical Information-Agnostic Flow Scheduling for Datacenter Networks,” in *HotNets 2014*.
- [15] W. Bai, L. Chen, K. Chen, D. Han, C. Tian, and H. Wang, “Information-Agnostic Flow Scheduling in Commodity Data Centers,” <http://sing.cse.ust.hk/papers/pias-tr.pdf>, Technical Report HKUST-CS15-01, 2015.
- [16] T. Benson, A. Akella, and D. A. Maltz, “Network Traffic Characteristics of Data Centers in the Wild,” in *IMC 2010*.
- [17] J. G. Carlsson and J. Shi, “A linear relaxation algorithm for solving the sum-of-linear-ratios problem with lower dimension,” *Operations Research Letters*, vol. 41, no. 4, pp. 381–389, 2013.
- [18] L. Chen, S. Hu, K. Chen, H. Wu, and D. Tsang, “Towards Minimal-Delay Deadline-Driven Data Center TCP,” in *HotNets 2013*.
- [19] A. Dixit, P. Prakash, Y. C. Hu, and R. R. Kompella, “On the impact of packet spraying in data center networks,” in *INFOCOM 2013*.
- [20] R. Fielding, J. Gettys, J. Mogul, H. Frystyk, L. Masinter, P. Leach, and T. Berners-Lee, “Hypertext transfer protocol—HTTP/1.1, 1999,” *RFC2616*, 2006.
- [21] A. Greenberg, J. R. Hamilton, N. Jain, S. Kandula, C. Kim, P. Lahiri, D. A. Maltz, P. Patel, and S. Sengupta, “VL2: a scalable and flexible data center network,” in *SIGCOMM 2009*.
- [22] C.-Y. Hong, M. Caesar, and P. Godfrey, “Finishing flows quickly with preemptive scheduling,” in *SIGCOMM 2012*.
- [23] K. Jang, J. Sherry, H. Ballani, and T. Moncaster, “Silo: Predictable Message Completion Time in the Cloud,” Tech. Rep. MSR-TR-2013-95, 2013.
- [24] V. Jeyakumar, M. Alizadeh, D. Mazieres, B. Prabhakar, C. Kim, and A. Greenberg, “EyeQ: practical network performance isolation at the edge,” in *NSDI 2013*.
- [25] B. Kalyanasundaram and K. R. Pruhs, “Minimizing flow time nonclairvoyantly,” *Journal of the ACM (JACM)*, vol. 50, no. 4, pp. 551–567, 2003.
- [26] A. Munir, G. Baig, S. M. Irteza, I. A. Qazi, A. Liu, and F. Dogar, “Friends, not Foes - Synthesizing Existing Transport Strategies for Data Center Networks,” in *SIGCOMM 2014*.
- [27] A. Munir, I. A. Qazi, Z. A. Uzmi, A. Mushtaq, S. N. Ismail, M. S. Iqbal, and B. Khan, “Minimizing flow completion times in data centers,” in *INFOCOM 2013*.
- [28] Y. Peng, K. Chen, G. Wang, W. Bai, Z. Ma, and L. Gu, “HadoopWatch: A First Step Towards Comprehensive Traffic Forecasting in Cloud Computing,” in *INFOCOM 2014*.
- [29] S. Radhakrishnan, Y. Geng, V. Jeyakumar, A. Kabbani, G. Porter, and A. Vahdat, “SENIC: scalable NIC for end-host rate limiting,” in *NSDI 2014*.
- [30] I. A. Rai, G. Urvoy-Keller, M. K. Vernon, and E. W. Biersack, “Performance Analysis of LAS-based Scheduling Disciplines in a Packet Switched Network,” in *SIGMETRICS 2004*.
- [31] K. Ramakrishnan, S. Floyd, and D. Black, “RFC 3168: The addition of explicit congestion notification (ECN) to IP,” 2001.
- [32] S. Schaible and J. Shi, “Fractional programming: the sum-of-ratios case,” *Optimization Methods and Software*, vol. 18, no. 2, pp. 219–229, 2003.
- [33] B. Vamanan, J. Hasan, and T. Vijaykumar, “Deadline-aware datacenter tcp (d2tcp),” in *SIGCOMM 2012*.
- [34] C. Wilson, H. Ballani, T. Karagiannis, and A. Rowtron, “Better never than late: Meeting deadlines in datacenter networks,” in *SIGCOMM 2011*.
- [35] H. Wu, J. Ju, G. Lu, C. Guo, Y. Xiong, and Y. Zhang, “Tuning ECN for data center networks,” in *CoNEXT 2012*.
- [36] Y. Xu, M. Bailey, B. Noble, and F. Jahanian, “Small is Better: Avoiding Latency Traps in Virtualized Data Centers,” in *SOCC 2013*.
- [37] M. Yu, A. G. Greenberg, D. A. Maltz, J. Rexford, L. Yuan, S. Kandula, and C. Kim, “Profiling Network Performance for Multi-tier Data Center Applications,” in *NSDI 2011*.