



**POLITECNICO  
DI TORINO**

**POLITECNICO DI TORINO**

Master Degree course in Communications and Computer Networks Engineering

Master Degree Thesis

# **Priority-driven load balance for low-latency datacenter networks**

**Supervisors**

Prof. Paolo GIACCONE

Prof. Andrea BIANCO

**Candidate**

Alessandro CORNACCHIA

Sessione di boh



# Acknowledgements

Thanks to my relatives and my loved ones for their unceasing support throughout my academic career. All of them had a fundamental role in giving me always new motivations to overcome step by step the tough obstacles I dealt with. Particularly, thanks to my parents who leave me the possibility of pursuing my studies. I hope to reward their sacrifices. Wide appreciation has to be granted to Politecnico di Torino and its personnel for the excellent teaching quality and the outstanding organization of its services. Finally, special acknowledgments go to German Sviridov, for his contributions of paramount importance.

## Abstract

Major cloud computing providers agree on denoting the flow completion time (FCT) as the primary goal to achieve in the design of a data center network (DCN). This is motivated by the fact that latency affects the performance of most cloud computer applications, such as web search, video streaming and social networking, with a direct impact on quality of experience as perceived by end users, hence ultimately on revenues.

Many existing approaches in literature rely on prioritization mechanisms at flow or packet level and schedulers in the network to reduce the average FCT, albeit for some applications also tail latency is important. To this purpose, the Shortest Remaining Processing Time (SRPT) scheduling policy is proven to be optimal when the job size is known in advance. Unfortunately this information is rarely available, therefore some proposals emulate the complementary scheme that give service first to flows that have transmitted less (Least Attained Service), exploiting a finite number of priority levels. All flows start with the highest priority and are progressively demoted to lower priorities as they receive service. The effectiveness of these algorithms, is largely improved augmenting the number of priorities. However, DCNs are usually realized with low cost commodity devices, where only few queues per port - typically two - are vacant.

The contribution of this work is to investigate the feasibility and evaluate the performance of a design which exploits multi-pathing - offered by DCN Fat-Tree topologies - to intelligently route flows across the switching fabric depending on their priority, so as to better segregate latency demanding flows. In other words, the key observation is that load balancing and prioritization can be jointly analyzed, in what multiple outgoing links together provide more priority queues than a port alone. Instead of assigning flow priorities with a scope limited to single interfaces, the problem can be addressed at data center level considering the set of links towards the flow destination as a whole. Unlike aforementioned techniques, longest flows are demoted across all priority queues of all exit interfaces, thus implying that routing choices are directly based on priority.

First, an analytical queueing model is provided for the setting of optimal parameters. Then, it will be shown through queueing model experiments that the proposed strategy indeed it is helpful under the few priority queue regime, but that long flows can be excessively penalized. Finally, extensive large scale packet level simulations in a real data center topology are conducted with Omnet++ discrete-event simulator, in order to validate the results obtained with previous steps.

# Contents

<b>List of Figures</b>	IV
<b>1 Datacenter networks</b>	1
1.1 Interconnection network design	2
1.2 Traffic engineering	3
1.2.1 Traffic control objective: FCT	3
1.2.2 Traffic properties	4
1.2.3 Challenges in traffic control	5
1.2.4 Datacenter TCP(DCTCP)	6
<b>2 Reducing FCT with flow prioritization</b>	7
2.1 Theoretical scheduling background	7
2.1.1 Flow-aware disciplines	7
2.1.2 Flow-agnostic disciplines	8
2.2 State-of-the-art solutions in practical networks	9
2.2.1 pFabric: targeting ideal schedulers	10
2.2.2 PIAS: the reference model	12
2.3 The problem of demotion thresholds	13
2.3.1 Stochastic queueing model	13
2.3.2 Greedy threshold assignment	16
<b>3 The spatial-diversity framework</b>	19
3.1 Improving scheduling with spatial-diversity	19
3.1.1 Abstraction as queueing system	20
3.2 Mathematical formulation for spatial-diversity	22
3.2.1 Decouple load balance from sub-thresholds	26
<b>4 Numerical analysis of spatial diversity</b>	29
4.1 Model implementation	29
4.1.1 Workloads	29
4.1.2 Optimal traffic load balancing	32
4.2 The service discipline: FIFO or PS?	37
4.2.1 Dimensioning spatial-diversity	37
4.2.2 Final comparison	41

<b>5</b>	<b>Evaluation in a Datacenter Network</b>	<b>43</b>
5.1	title . . . . .	43
	5.1.1 title . . . . .	43
	5.1.2 title . . . . .	43
5.2	title . . . . .	43
	5.2.1 title . . . . .	43
	5.2.2 title . . . . .	43
	<b>Bibliography</b>	<b>45</b>

# List of Figures

1.1	interconnection networks. . . . .	2
2.1	Datacenter fabric abstraction as a giant switch . . . . .	10
2.2	PIAS overview . . . . .	12
2.3	Queueing model . . . . .	14
3.1	Demotion extension with spatial diversity. Red arrows indicate demotion trajectory of the longest flow. On each switch is reported the priorities it handles. Lower indexes (and darker colors) correspond to higher priorities. . . . .	20
3.2	Three queueing system comparison. . . . .	21
3.3	Mapping of the model abstraction with DCN topology . . . . .	21
3.4	Job size distributions observed by servers $s_1, s_2, s_3$ in decreasing order of priority. Optimal load balance split for $\lambda = 0.9$ . . . . .	24
4.1	Workload properties . . . . .	30
4.2	Web search workload. Simple case of $K=2$ . . . . .	33
4.3	Data mining workload. Simple case of $K=2$ . . . . .	33
4.4	M/M/1 response time . . . . .	34
4.5	Web Search workload . . . . .	35
4.6	Data Mining workload . . . . .	36
4.7	Web Search workload and FIFO discipline at 99% confidence interval. . .	38
4.8	Web Search workload and PS discipline at 99% confidence interval. <b>GRAFICI MANCANTI MA SIMULAZIONI GERMAN CON QUELLO SCENARIO GIA FATTE MA USATE NEI GRAFICI FINALI. CHIEDERE GRAFICO A PARTE SE UTILE</b> . . . . .	39
4.9	Web search workload and PS discipline at 99% confidence interval . . . .	40

# Chapter 1

## Datacenter networks

In the last decades, with the spreading of cloud services accessible to anyone, computing has undergone a remarkable evolution. Encouraged by the astonishing growth of virtualization technologies in the IT industry, as well as the availability of storage and chips at ever-more modicum prices, over-the-top (OTT) players like Google, Amazon and Microsoft, have been building datacenter hotspots all around the world. Most of the applications any sector of modern society relies on, such as commercial and financial services, Web search, scientific computing, on-demand video streaming, recommendation systems, not to mention social networking and online gaming, more often than not run inside one of their datacenters' infrastructure. Indeed, from small enterprises to large corporations, it has been a common cost-saving strategy to offload, up to a certain extent, the deployment and operation of their own information systems to third parties, the cloud service providers. After all, it is well-known in engineering, how much resource concentration enables better design and ease of maintenance.

Among the major challenges posed by cloud computing paradigm, certainly there is the design of a huge communication network, hosting hundreds of thousands of servers, that is required to simultaneously provide high throughput and low-latency, while guaranteeing uninterruptible service continuity and lending itself to relentless expansion. The peculiar requirements of the datacenter environment and the lack of flexibility of traditional TCP/IP stack, have pushed the development and adoption of unconventional approaches, such as centralized control and network softwarization (SDN), laying the premises for a wider transformation process in the network industry. In this sense, datacenters can be credited to having represented natural incubators for the evolution of telecommunication networks happening during the last few years. As a matter of fact, virtualization technologies and control plane programmability soon would have become disruptive innovations in networking ecosystem, that started to experience the same changes the IT world faced years before. As a result of this process, nowadays, ISPs are converting their infrastructure towards the same solutions, re-architecting PoPs as small-scaled modern datacenter, with massive employment of virtualization as regards network functions and devices (NFV). Similarly, the next generation mobile network, 5G, is going to base both its core and edge functions on the same paradigm, as revealed by its standardization and by the investments in multi-access edge computing (MEC).



## 1.1 Interconnection network design

The typical structure of a datacenter, as shown in Fig.??, is comprised of many *racks*, interconnected among each other thanks to a common network infrastructure. A rack is nothing else than a group of servers physically co-located in a common cabinet, attached to the same *Top of Rack* switch (ToR) and thus separated by a single hop. In achieving the goal of hosting an enormous amount of servers, the principal bottleneck often results in being the interconnection fabric, usually referred to as *Data Center Network* (DCN). Ideally, it should act as a huge switch, able to provide maximum-rate communication among servers, that is their NIC's access capacity. At high level, there are fundamentally two ways for practically realizing such infrastructure. The first choice is to rely on complex specialized solutions, like InfiniBand, or high-performance IP devices with many ports, that successfully provide bandwidth for thousands of nodes, however incurring in high deployment and management costs. Conversely, the second possibility is to build the network infrastructure by simply leveraging on commodity off-the-shelf switches that are cheaper, already on economy of scale and fully compatible with existing hardware and operating systems, just as large distributed clusters are made of general purpose cheap computers. This is usually the design pursued by principal cloud providers, as recently disclosed by themselves. A major drawback of the latter strategy, is the difficulty in providing full-rate communication among hosts in different racks. In other words, depending on traffic patterns and especially network topology, it is very hard to place enough inter-rack capacity - usually referred to as *bisection bandwidth* - to satisfy the collective demand of all the racks. For this reason, the topology design plays a fundamental role for the feasibility of scalable and cost-beneficial large datacenters. It is worth noticing that sometimes, building a DCN with full bisection bandwidth is unnecessary and cost prohibitive, therefore it is commonly accepted some degree of *oversubscription*, meaning that the ratio between total intra-rack and bisection bandwidth is greater than 1. This choice is acceptable when applications are mostly rack local and statistical multiplexing is beneficial. In general, a non-oversubscribed network allows greater flexibility in resource allocation and deployment of applications across racks, leading to higher utilization which is also important for big cloud providers.

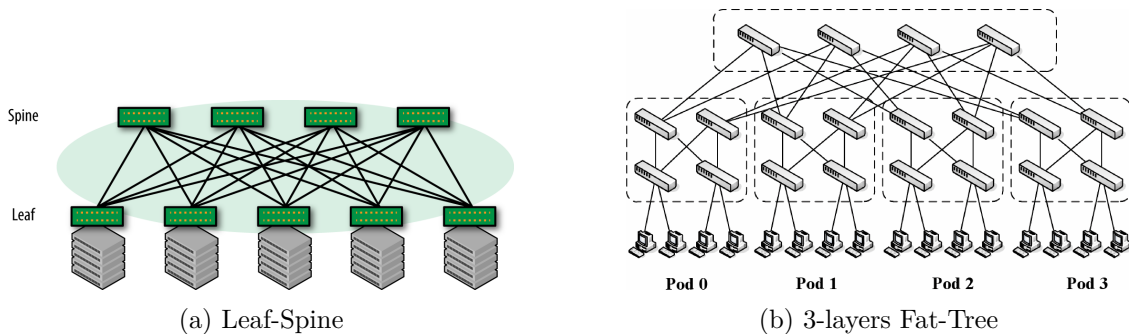


Figure 1.1: interconnection networks.

Two examples of topology design are provided in Fig. 1.1. The *Leaf-Spine* (Fig. 1.1a),

is a typical hierarchical architecture, where the interconnection between leaf and spines switches form a bipartite graph. Despite being popular in campus networks, scaling to thousands of servers would require having spines with lots of high-capacity ports. A more flexible solution dates back to the early telephone network, when Charles Clos had to solve a similar problem and came up with its proposal for multi-stage switching fabrics. Indeed, the *Fat-Tree* (Fig. 1.1b), representing the most widely adopted topology, is a folded Clos network. The main plus of this configuration is the possibility to host as many servers as desired, using only switches with a fixed number of port  $k$ , eventually providing full bisection bandwidth, just by recursively adding new layers, or stages. Fig. 1.1b shows an example of 3-layer Fat-Tree with 4 ports switches. Essentially, it is a recursive design, in which a  $l$ -layer network is built by connecting  $k$  blocks, called *Point of Delivery* (POD). Note that in terms of PODs, a Fat-Tree is actually a Leaf-Spine with  $k/2$  spines and  $k$  POD leafs. Each POD has the same structure of a  $(l-1)$ -layer network, unless for the fact that  $k/2$  ports from those nodes that were the spines of layer  $l-1$  have been used for interconnecting the POD to the new spines of the  $l$ -layer architecture. To put it differently, the  $l$ -layer DCN can become a POD for the  $(l+1)$ -layer DCN by disconnecting  $k/2$  PODs from the spines to and connecting them to a new row of spines (core switches). The elementary POD is the 2-layer architecture, that can host at most  $k/2)^2$  servers. Thus, a datacenter with  $l$  stages can support up to  $k^l/2^{l-1}$  machines, and this is the reason behind popularity of Fat-Tree, other than its compatibility with main requirements for datacenters, such as path diversity, which is important for fault-tolerance and traffic engineering. Throughout this work, in particular, it will be investigated a flow scheduling solution that exploits this multi-path property.

## 1.2 Traffic engineering

Large data centers represent a challenging environment for network designers. They host many thousands of servers running a myriad of applications belonging to tenants with heterogeneous QoS demands on a very physically circumscribed network. Data travel for no more than few hundreds meters with negligible propagation delays, through links of huge capacities up to 40Gbps. In such a context, the traditional TCP/IP protocol stack alone operates inefficiently, hence custom transport protocols, traffic control and traffic management techniques have been deployed. In this section will be first pointed out common objectives for network operations, then briefly reviewed the traffic characteristics and the subsequent issues they pose, in order to lay down the fundamentals of this work and justify the choices for traffic generation undertaken in its simulative part.

### 1.2.1 Traffic control objective: FCT

There is unanimous consensus among data center providers in putting effort to minimize the *Flow Completion Time* (FCT) metric, that is the all-up delay since when data are requested by a client application to the time they are at its disposal. The reason behind this interest is that latency directly impacts the quality of experience (QoE) as perceived by end users and ultimately operator revenues. For interactive Web applications and online services, the responsiveness determines the number of users on the long run. Also,

low-latency communication also enables flexible inter-rack deployment of micro-services across the network. Therefore, a common goal is to find efficient traffic control algorithms to handle latency sensitive flows, while maximizing network utilization which is also important for cost savings. Usually, different techniques are compared by measuring the *average* flow completion time.

### 1.2.2 Traffic properties

Several studies have been conducted in literature to characterize the main properties of traffic in data centers. Traffic characteristics are highly dependent on applications, that determine flow sizes, flow arrival patterns and the requirements from a network perspective. Common applications running in datacenters and for which there exist traffic studies are Web search, data mining (e.g. Hadoop) and cache services. For Web search queries and the corresponding responses, for example, few packets are enough, so they generally comprise short flows. Instead, other services such as data mining and batch computing tasks, may transfer large amount of data. Additionally, long-lived background connections of large size are continuously present for VM migration, backups, consistency updates and data replication. A common scenario observed by different studies is that the majority of flows are short, but overall they do not significantly contribute to the total traffic, which is mainly carried by few large flows. Measurements from a large cloud provider production datacenter reveal that 80% of flows are less than 10KB and almost all flows (99%) are less than 100MB. However, more than 90% of transferred bytes are in the 1% of flows greater than 100MB. Facebook [12] shared its traffic statistics and reported the median flow size to be 3KB for Web Search and 100KB for data mining with Hadoop, while the tail flow size 10MB and 100MB respectively. Similar trends are claimed by other authors [banson, pfabric, dctcp], although flow sizes are slightly different depending on datacenter services.

Therefore, the first remark is that a variety of flows with different sizes coexists on the same DCN, the majority of which lasts a couple of packets only.

Frequently, different flow sizes happen to be associated with different QoS requirements, whose knowledge drives in the design of transport protocols and traffic control algorithms and meet application demands. Short flows are typically query traffic, very sensitive to latency, stemming from the *Partition/Aggregate* pattern which is widespread in distributed computing, from social networking content composition to retail and recommendations. According to this paradigm, a task requested at higher layer to an *Aggregator*, is broken into simpler units that are dispatched along a tree-like logic to lower level aggregators, that may further break the request into smaller pieces, until they're finally handled by worker nodes. The responses of the workers, when ready, are then conveyed back along the same reversed tree logic to aggregators that put together the results. Key in this process is that it must complete within strict deadlines, on the order of 10-100ms, that are determined in order to satisfy the worst-case latency tolerated by the Service Level Agreement (SLA) with customers, or tenants in cloud computing jargon. Ideally, application developers should not be concerned of network delays and should be entitled to employ the time before deadlines to improve final results thus end

user satisfaction, without resorting to the implementation of complex ad-hoc solutions to compensate for network inefficiency. On the other hand, long flows are mainly comprised of update flows that carry fresh data for parallel computing jobs (e.g. MapReduce), or transfers for data replication across servers, also located in different facilities in the globe. They are throughput-oriented flows, demanding considerable bandwidth, but they are not sensitive to delays. Section ?? will highlight some issues deriving from the mix latency critical flows arising from the partition/aggregate workload with background long-lived connections.

Finally, one may also be interested in understanding the communication patterns to tailor traffic engineering choices and perform capacity planning. For instance, knowing the degree of traffic rack locality would allow better awareness in deciding the oversubscription ratio. To this extent, it is difficult to draw general conclusions and prior studies have given contradictory results: some work in literature ( cite traffic-in-the-wild ) reports a marked locality, whereas some other observes completely different patterns with traffic not at all rack local. The reason lies in how applications are deployed across servers and clusters. In the Facebook data center each machine is assigned a precise role and machines with a same role are grouped in the same rack. Since Web Servers machines talk primarily to Cache Followers machines, there is substantial intra-rack communication for this service. Conversely, Hadoop traffic stays 75% of the times in the same rack. What is true in general is that flow arrival rates are quite high, as many as thousands flows per second per server. Combined with the fact that the majority of them is short and that their destinations is often randomized at application level to avoid hot-spots, it follows that the traffic matrix of a datacenter network has been revealed to be very fluctuating, unstable and difficult to predict.

### 1.2.3 Challenges in traffic control

The traffic characteristics illustrated so far, produce undesired pattern and effects on data center networks, which standard TCP transport suffers specifically. They are *queue buildup* and *incast* problems.

1. *Queue buildup*. On the basis of well-known congestion control mechanisms, traditional TCP sources increase their window until they experience either a timeout or a packet drop, resulting in the usual sawtooth pattern. This causes switch queues to grow, mainly due to long connections which have time to inflate enough their window. In the context of traffic depicted in the previous section, this is especially problematic because long flows sharing the same queues with short flows harvest much of the buffer space, causing short flows either to queue behind them and increase their latency, or to experience more penalizing packet drops. The queue buildup impairment is even more severe in data center networks, since commodity switches are generally cheap shared buffer architectures, where packets from different ports are stored logically in the same memory space. Thus, high utilization of a single port could degrade the performance of flows traversing a different interface of the same device.
2. *Incast*. The incast problem arises from the Partition/Aggregate pattern common

in data centers applications. After the aggregators assign portions of the same task to different workers, the flows containing the responses tend to be cluster at the same time on the same switch ports, on the way back towards aggregators. Such a flow synchronization, in conjunction with the queue buildup phenomenon, causes synchronized drops even if the flows are short, as in the case of responses to Web queries. Packet drops, therefore timeouts, are not acceptable for deadline constrained flows, typical of the Partition/Aggregate pattern.

#### 1.2.4 Datacenter TCP(DCTCP)

Many techniques for traffic control are possible in order to counteract the previous impairments: they include transmission rate control, traffic shaping, routing and load balancing, scheduling and priority schemes. In the following of this work it will be explored a solution which combines prioritization and routing. However, a milestone for transport protocol design in data centers has been Datacenter TCP (DCTCP). It actually consists in a very concise modification to TCP which leverages Explicit Congestion Notifications (ECN) from the network to properly modulate the window size of the sources.

##### DCTCP

The basic idea behind DCTCP is that queues in the network should be kept as empty as possible to avoid large backlogs that increase latency and do not leave enough headroom to absorb burst arrivals, occurring for example due to incast. Instead of pushing the window to grow until a packet drop is detected, the DCTCP transmitter slows down proactively depending on the level of congestion on the bottleneck link. Network switches mark every packet with a congestion signal as soon as the queue where the packet is stored exceeds a given occupancy  $K$ . The TCP receiver then conveys back the congestion signaling to the TCP transmitter, setting the ECN-Echo bit to 1 in its next ACKs. Finally, the TCP sender maintains an estimate  $\alpha$  of the fraction of marked packet on an interval of roughly one RTT and modulates the window as:

$$cwnd = cwnd \times (1 - \alpha/2)$$

This way the window size is gently reduced upon mild congestion — note that only in case  $\alpha=1$  it is cut in half as in standard TCP — still ensuring high throughput, but reducing its aggressiveness.

It has been proven that DCTCP effectively succeed in lowering the amplitude of queue oscillations to  $O(\sqrt{BDP})$  from the  $O(BDP)$  of TCP, being BDP the bandwidth-delay product, while not losing throughput for a proper setting of the marking threshold  $K$ . Note that the only requirement from the network is to configure switches with an AQM scheme to mark packets. In practice this can be achieved configuring RED, already available in most devices, so that it marks based on the instantaneous queue length and with a unique high and low threshold equal to  $K$ .

## Chapter 2

# Reducing FCT with flow prioritization

To the purpose of minimizing the average flow completion time (FCT), a common strategy is to treat short flows with tight latency constraints differently from the others. This can be managed with prioritization and scheduling algorithms, for which there already exist numerous theoretical studies. In the following of this chapter will be reviewed the most significant of them, with care to their application to datacenter networks.

### 2.1 Theoretical scheduling background

At high level, scheduling policies can be classified into two categories, depending whether or not flow properties — such as size and deadline — are known a priori.

#### 2.1.1 Flow-aware disciplines

Flow-aware disciplines are those techniques that assumes the job characteristics are precisely known before initiation. While this is often the case for many systems in other industries like manufacturing, it is not always so in communication networks. Sometimes the flow size is known exactly, for example when a server is requested a static object (e.g. a static Web page, a file transfer) or it can be roughly estimated, but generally speaking this information is either not available or it involves undesired modifications of the application layer.

The optimal approach to minimize the job completion time for an offline system is the *Shortest Job First* (SJF) discipline that consists in serving jobs in decreasing order with respect to their size. However this policy is not suitable for dynamic contexts where new jobs can arrive at any time instant. For such scenarios has been adopted a preemptive version of SJF, known as *Shortest Remaining Processing Time* (SRPT), which chooses first the job with shortest time left to its completion, or equivalently in the context of flow scheduling, the smallest amount of bytes left. SRPT has been proven for long to be the optimal policy for minimizing the average response time (i.e. FCT) in a single server system, regardless of the serving time and inter-arrival distribution [13].

### 2.1.2 Flow-agnostic disciplines

In absence of precise knowledge about the flow length, a very effective policy is the dual approach to SRPT, the so called *Least Attained Service* (LAS) or *Foreground-Background* (FB). LAS is a preemptive scheduler which gives service first to the flow that has transmitted less bytes, serving in processor sharing when there are ties among flows. In other words, a job retains alone the processor until it has received the same service of another jobs in the system or it is preempted by a newly arrived shorter job. If no fresh jobs arrive, those in the systems prosecute sharing the processor. The main insight of LAS is to exploit the increasing knowledge about the flow size gained during its service. In fact, the LAS scheduler becomes more and more confident that a given flow is a large one, as further of its bytes are transmitted.

LAS is optimal among the flow-agnostic disciplines when the hazard rate of the flow size distribution is a decreasing function [1]. The hazard rate is defined as the ratio of the probability density function  $f_X(x)$  to the survival function:

$$h(x) = \frac{f_X(x)}{1 - F_X(x)}$$

and it is a function that represents the instantaneous failure rate of a quantity. For instance, it can be seen as the likelihood that the flow size  $X$  ends at value  $x$  given that  $x$  bytes have already been observed. As a general rule, heavy-tailed distributions exhibit a decreasing  $h(x)$ , whereas  $h(x)$  is increasing for light-tailed distributions and constant for the exponential distribution. Thus, LAS works at best under job size distributions that present high variability, which is a case of interest since they well model datacenter's traffic where a majority of short flows coexists with few very large flows.

At this point, it is interesting to better understand how LAS compares to other disciplines, such as PS, not only on average but in depth with respect to flow sizes, in order to quantify its impact on small and large jobs. Also, it would be useful to bound the sub-optimality of LAS with respect to the SRPT flow-aware scheduler. To this purpose, the authors of [11] compared different distributions with increasing coefficients of variation, under LAS, PS and SRPT disciplines. The *coefficient of variation*  $C$ , is a single number measure of the variability of a distribution and it is expressed as the standard deviation  $\sigma$  normalized to the mean of the distribution  $\mu$ :

$$C = \frac{\sigma}{\mu}$$

The comparison was based on the *average conditional completion time*, that is the average completion time of flows belonging to a given size. Formally, denoting as  $T$  the random variable associated to the completion time and  $X$  as the random variable associated to the flow size, the average conditional completion time is defined as  $\mathbb{E}[T|X = x]$ . This is a rather practical quantity to tell how the FCT vary with flow sizes and to highlight unfairnesses brought by the scheduling policy.

The distributions taken into account were the negative exponential distribution, whose probability density function rapidly converges to zero, and a set of Bounded Pareto distributions with varying  $C$ , as example of heavy-tailed distributions. Their probability density functions are given by:

$$\begin{aligned}
f(x)_{Exp} &= u(x)\mu e^{-\mu x}, & \mu &\geq 0 \\
f(x)_{BP} &= \frac{\alpha k^\alpha}{1 - (k/p)^\alpha} x^{-(\alpha+1)}, & k \leq x \leq p, 0 \leq \alpha \leq 2
\end{aligned}$$

For the Pareto distributions, different setting of their parameters correspond to different  $C$  values, but generally it holds  $C \gg 1$  which means essentially high variability. For the negative exponential distribution  $C$  is constant and equal to 1.

- **LAS vs PS.** They found that for Pareto distribution with  $C \geq 6$ , LAS outperforms PS above the 99th percentile of the flow size distribution, that is the conditional completion time  $\mathbb{E}[T|X = x]$  is better under LAS than under PS for all but the largest 0.3% flows. Additionally, even the penalty of the longest flows is within a factor 2. Notice that it is reasonable to expect some slowdown for very long flows due to starvation, especially when an elephant flow meet a longer one, which is queued behind it. Instead, for exponential distribution, flows with size above the 80th percentile are severely penalized using LAS, but overall the average FCT is still better.

The final message to be derived here is that LAS is always beneficial for the average flow completion time in respect of PS for distributions with  $C \geq 1$ .

- **LAS vs SRPT.** This comparison is of interest because it shows quantitatively the performance gap with the optimal policy. Indeed, it was provided an analytical expression for the worst case penalty of LAS with respect to SRPT in terms of mean conditional completion time, valid for every continuous time finite mean and finite variance distribution. In particular, when applied to Pareto, LAS is very close to SRPT for all job sizes, with a penalty of 1.25 under all load conditions. This result is of remarkable importance as essentially states that for heavy-tailed distributions with high variability there is no significant gain in knowing the flow size beforehand, but good performances can be achieved with the LAS policy, which is simpler to implement.

In short, the ideal scheduler is SRPT if detailed flow information are available at transport upon flow initiation, alternatively LAS is the best choice provided that flow sizes present high variability. Any practical design proposed in literature and revised in the following sections aim to approximate any of these targets.

## 2.2 State-of-the-art solutions in practical networks

Theoretical results, both for flow-aware and flow-agnostic scheduling disciplines, refer to a scenario with a single link and do not account at all for the implementation complexity. A major restraint of LAS and SRPT that limits their practical applicability are the fine-grained decisions they adopt for flow scheduling. To determine the next job to serve, it is required to maintain per-flow state information and perform comparisons among them



at every packet transmission. These operations can be very complicated to implement at high line rates with thousands of simultaneous flows, as typically happens in DCNs. (Sec. ??). Moreover, in a data center network many servers are connected through a multi-stage switching fabric (1.1), which implies that, actually, from source to destination multiple links are traversed, hence optimizing local decisions does not ensure global optimality. As a trivial example, consider a situation with three server  $s_1$ ,  $s_2$  and  $s_3$  and three flows  $f_1 : s_1 \rightarrow s_2$ ,  $f_2 : s_1 \rightarrow s_3$  and  $f_3 : s_2 \rightarrow s_3$ . With this setting  $f_1$  competes against  $f_2$  in the access interface of  $s_1$ , while  $f_2$  competes with  $f_3$  in the egress towards  $s_3$ . If a local scheduler in the outgoing interface of  $s_1$  decide to serve  $f_2$  before  $f_1$ , but a second independent scheduler in the last interface ahead of  $s_3$  prioritize  $f_3$  over  $f_2$ , the choice of the former scheduler would be vanished by the conflicting choice of the latter, because  $f_2$  would delay  $f_1$  despite being queued back  $f_3$ . Therefore, the optimal scheduling pattern could be achieved by means of a global network view.

The state-of-the-art proposals that address these issues are pFabric [3] and PIAS [4]. They have been taken as the reference model of this work, which will try to extend their fundamental insights.

### 2.2.1 pFabric: targeting ideal schedulers

Alizadeh et al. provide a general representation about the problem of scheduling flows over the datacenter fabric. Specifically, they abstract the DCN as a giant switch inter-connecting all the servers, as shown in Fig. 2.1. In this representation all the links are assumed to be unidirectional, so that leftward there are the servers' access NICs and rightward the ToR egress interfaces towards the same servers.

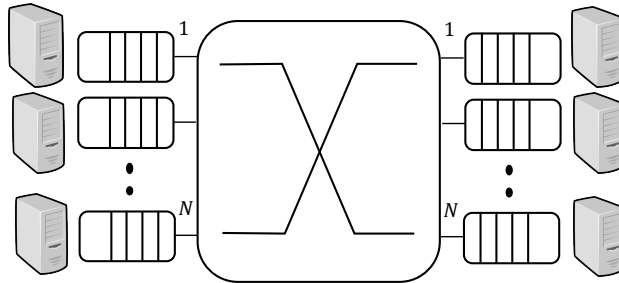


Figure 2.1: Datacenter fabric abstraction as a giant switch

### Ideal flow scheduling in DCN

Assuming that the fabric can sustain maximum throughput and that flows compete against each other only at ingress and egress interfaces, the optimal solution to minimize the flow completion time can be found solving a NP-hard problem known as *sum-multicoloring*. Nevertheless, a simpler greedy algorithm exists and it has been proven to be very close to the optimum. Authors referred to this algorithm as the *Ideal* algorithm and used it as a baseline for the evaluation of pFabric. Briefly, it consists in a maximal flow scheduling, where at every new flow arrival or depart, flows are sorted in ascending

order of data remaining up to their completion and served with this order. Flows are not served until there is another flow with less data remaining traversing the same ingress or egress port. Despite being a simplified model, still the *Ideal* algorithm is a valid benchmark for any design that targets flow completion time minimization, since it is a sort of best case latency at a desired level of load and over an ideal interconnection network. At this point, pFabric showcased that it is possible to achieve nearly ideal performances in a distributed way, treating each interface on its own with local scheduling decisions and simultaneously rely on minimal transport strategies. In particular, their work can be summarized in the following key aspects.

1. *Knowledge of flow requirements.* It is assumed to know at transport layer flow sizes or flow deadlines for deadline-constrained traffic.
2. *Prioritization.* Schedulers in the network are priority schedulers. Depending on the scheduling objective, the packets encode with a priority a different metric, on which the scheduler choices are based on. For example, to approximate SRPT on every single link, the priority may indicate the amount of bytes still to transmit for a given flow. Packets are dequeued and dropped according to their priority.
3. *Simple rate-control.* Rate control at end hosts is minimal and aims just to avoid persistent congestion. In fact, contrary to DCTCP (Sec. 1.2.4), even if queue sizes grow, only the performance of few long flows is significantly impacted, since it is used a detailed prioritization mechanism both for packet transmission and drops. Some shrinkage to standard TCP have been described to realize such minimalistic transport.

### Approximate optimality with priority queues

As already mentioned, implementing such a prioritization scheme on available commodity devices is very challenging and for sure would require hardware modifications. The interested reader could find a possible digital design in the paper of pFabric. The same source, however, provided also a straightforward solution readily deployable with current switches. The idea is to coarse the granularity of the scheduler by adopting only a finite number of priority levels. Then, packets with different priorities are enqueued in corresponding priority queues (PQs) and handled with traditional network schedulers (SP, WRR, WFQ, etc.). In pFabric it is employed a Strict Priority (SP) scheduler, but in general other choices are not precluded [5]. In short, LAS and SRPT schedulers, as well as other disciplines, can be emulated by tagging packets with a priority label and leveraging separated queues in network interfaces. The typical number of priority queues available in datacenter switches ranges between 2 and 8, albeit it often occurs that some of them are reserved — more details in Sec. 2.2.2. Of course, increasing the number of PQs results in better approximation of the ideal scheduling, that implicitly would correspond to having an infinite number of PQs. In fact, the ideal scheduler directly compares each other all priorities of the packets in the buffer.

Finally, the FCT gain obtained with this system largely depends on the way flows are clustered in priority levels. This underlies the need of a careful tuning of a set of thresholds which split packets among the priorities available. A dedicated section of this work illustrates some criterion to choose them (§??).

### 2.2.2 PIAS: the reference model

The next step with respect to pFabric as well as the starting point of this work is PIAS: Practical Information-Agnostic Scheduler [4]. This proposal is the first which addresses the scheduling problem without the assumption of an exact prior knowledge about the flow size. Instead, PIAS tries to mimic a LAS scheduler exploiting the sole knowledge of the distribution of flow sizes rather than their precise values. Surprisingly, when compared to pFabric it delivers very similar performances for short flows, that are the more critical ones (the gap is within 4.9%). At the same time, however, PIAS preserves ease of deployment with the current switch hardware and the standard distributed congestion-control algorithms. Summarizing, the two main design principles of PIAS are:

1. *Flow-agnostic.* It requires only the knowledge of DC-wide flow size distribution, which can be easily estimated once and updated dynamically. Authors do not account for any heterogeneity of the distribution across different racks [12], in what the problem would become quite difficult to treat with analytical methods. Also dynamic changes of the distribution along time have been ignored, but the architecture is flexible enough to adapt to this case.
2. *Low complexity.* It should be compatible with legacy TCP/IP protocol stack and readily deployable without touching the hardware of existing devices.

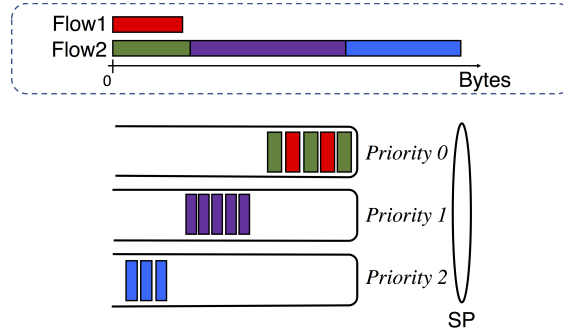


Figure 2.2: PIAS overview

PIAS embraces a *Multi Level Feedback Queue* (MLFQ) mechanism to resemble the LAS policy. MLFQ essentially apportions flows in a finite number of priority queues as proposed in pFabric, but in absence of flow size information flows are dynamically moved across priority queues. In more details, each packet of a given flow is mapped to a priority level, inversely proportional to the bytes it has already sent since its beginning. All flows start with the highest priority, then longest flows are progressively demoted to lower

priorities. Packets of same priority are buffered in the corresponding PQ according to a First-In First-Out (FIFO) order, while packets belonging to different priority queues are scheduled in Strict Priority (SP). In this way, packets belonging to short flows are always prioritized over those belonging to long flows, mimicking LAS but without incurring in the complexity of its implementation, which would require many comparisons. Figure 2.2 clarifies the demotion mechanism.  $Flow_1$  is a mice flow and it is entirely served at the highest priority, whereas  $Flow_2$  is an elephant flow and it is gradually de-prioritized up to the last PQ. Notice that differently from pFabric, all flows traverse highest priorities in their initial lifetime, since their size is initially obscure. Hence, longer jobs constitute a small impairment for short ones, that motivates the gap with pFabric.

A key point for this kind of systems is how to choose the demotion thresholds. In fact, especially when few queues are available, an unbalanced threshold setting may lead to severe performance degradation. On one hand, thresholds too small cause premature flow demotion and delays for short flows that get mixed with long flows, on the other hand if the thresholds are too large, medium and elephant flows overstay in high priorities, resulting again in worse FCT for delay sensitive flows. The next section will present simple heuristics and a more refined optimization to find the set of thresholds.

## 2.3 The problem of demotion thresholds

The same authors of PIAS proposed a queueing model to mathematically describe the dynamics of the system. It has to be remarked that the queueing model captures only the average flow completion time on a single interface. Therefore, it is assumed that the flow size distribution is homogeneous over the datacenter fabric so that bottleneck links observe all the same distribution. With this assumption, the average FCT over the whole fabric is just a linear rescaling of the average FCT on a single link, therefore the performance index of any set of thresholds obtained with the model is meaningful for the whole DCN as well.

In the following of this section, first it is formalized the queueing model along with its parameters, then it is formulated a non-linear minimization problem that can be used to optimize demotion thresholds. Finally, two other trivial heuristics for threshold assignment are reviewed.

### 2.3.1 Stochastic queueing model

The system, shown in Figure 2.3, is thought as a tandem of  $N$  subsequent M/M/1 queues. The customers are flows of size  $X$  extracted from a given distribution with cumulative function  $F(x)$  and arriving according to a Poisson process of intensity  $\lambda$ . Queues are lazy and able to serve at most a maximum amount of bytes for each flow, that depends on the demotion thresholds. When customers enter the system they are initially served by the first queue, then either they leave the system if all of their bytes have been served, or they are demoted to subsequent queue, up to queue  $N$ .

Let be  $Q_p (1 \leq p \leq N)$  the  $N$  priority queues. Denote with  $\alpha_i$  the threshold after which a flow changes its priority from  $i - 1$  to  $i$ , where higher priorities correspond to lower indexes and  $i \in [0, N]$ . The upper threshold is always set to  $\alpha_N = \infty$  so that all the

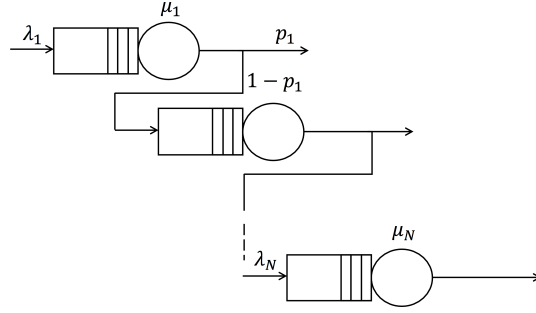


Figure 2.3: Queueing model

largest flows stay in the last queue, similarly  $\alpha_0 = 0$  for simplicity. Each queue  $Q_i$  can serve at most  $\alpha_i - \alpha_{i-1}$  bytes of any customer. In the real MLFQ system (Sec. 2.2.2), all queues of a given port are orchestrated by a Strict Priority (SP) scheduler, however this would make the queues dependent on each other complicating the analysis. To avoid the issue, the queues are treated as if they are independent and their priority hierarchy is taken into account by adjusting their drain rate  $\mu_i$ . In fact, if  $\mu$  is the overall link capacity, subsequent queues in the tandem are assigned a drain rate  $\mu_i$  equal to the residual bandwidth left after servicing previous queues. In practice, denote with  $\rho_i$  the load insisting on  $Q_i$ , then

$$\mu_i = \mu \prod_{j=1}^{i-1} (1 - \rho_j) \quad (2.1)$$

Next, for the sake of conciseness let's indicate with  $\theta_i \triangleq F(\alpha_i) - F(\alpha_{i-1})$  the probability that a new flow has size in  $[\alpha_{i-1}, \alpha_i)$  and abbreviate with the notation  $\bar{F}(x) = 1 - F(x)$  the survival function of  $X$ . The target is to derive the subsequent arrival rates  $\lambda_i$ . After service in queue  $i$ , a customer leaves the system with a probability  $p_i$ , that depends on flow size distribution. In fact, customers who leaves the systems after the  $i$ -th queue are the flows with size in  $[\alpha_{i-1}, \alpha_i)$ , among those with size in  $[\alpha_{i-1}, \infty)$ . Therefore,  $p_i$  can be expressed renormalizing the probability  $\theta_i$  as:

$$p_i = \theta_i / \bar{F}(\alpha_{i-1})$$

Consequently, from the definition of  $\theta_i$

$$1 - p_i = \bar{F}(\alpha_i) / \bar{F}(\alpha_{i-1})$$

If the system is stable, the arrival rate of queue  $i$  is nothing else than the output rate of queue  $i - 1$ , weighted by the probability of remaining in the system. Therefore:

$$\lambda_i = \lambda_{i-1} (1 - p_{i-1})$$

By iterative substitution it is possible to get a general expression for every  $\lambda_i$  that depends only on the flow generation intensity  $\lambda$  and the arrival rate at the previous queue  $\lambda_{i-1}$ .

Indeed:

$$\begin{array}{ll}
 \lambda_1 = \lambda & 1 - p_1 = \bar{F}(\alpha_1)/\bar{F}(\alpha_0) = \bar{F}(\alpha_1) \\
 \lambda_2 = \lambda_1(1 - p_1) = \lambda\bar{F}(\alpha_1) & 1 - p_2 = \bar{F}(\alpha_2)/\bar{F}(\alpha_1) \\
 \lambda_3 = \lambda_2(1 - p_2) = \lambda\bar{F}(\alpha_2) & 1 - p_3 = \bar{F}(\alpha_3)/\bar{F}(\alpha_2) \\
 \dots & \dots
 \end{array}$$

In general:

$$\lambda_i = \lambda\bar{F}(\alpha_{i-1})$$

Finally, this rate refers to flow arrivals, whereas thresholds  $\alpha_i$  are expressed either in bytes or in packets, because they serve for demotion. It is simple to obtain the byte arrival rate by scaling  $\lambda_i$  (flows/sec) of the average traffic size produced by these flows on queue  $i$ . Remember that due to the demotion mechanisms, customers of each queue are truncated versions of the original flows, whose size ranges in  $(0, \alpha_i - \alpha_{i-1})$ . and that, subsequent priority queues observe a flow distribution truncated above  $\alpha_{i-1}$ .

It is needed to compute  $\mathbb{E}[L_i]$ , the average length of customers served in queue  $i$ . It holds:

$$\mathbb{E}[L_i] = \underbrace{\int_{\alpha_{i-1}}^{\alpha_i} (x - \alpha_{i-1})f(x)dx}_{(i)} + \underbrace{(\alpha_i - \alpha_{i-1}) \int_{\alpha_i}^{\infty} f(x)dx}_{(ii)} \quad (2.2)$$

(i) Traffic generated by flows with sizes in  $[\alpha_{i-1}, \alpha_i]$

(ii) Traffic generated by flows larger than  $\alpha_i$

Define the truncated probability density function seen by queue  $i$  as:

$$f_i(x) = f(x) / \bar{F}(\alpha_{i-1})$$

It holds:

$$\lambda_i = \lambda \bar{F}(\alpha_{i-1}) \frac{\mathbb{E}[L_i]}{\bar{F}(\alpha_{i-1})} = \lambda \mathbb{E}[L_i] \quad (2.3)$$

Summarizing, it was possible to write down all arrival rates  $\lambda_i$  and all draining rate  $\mu_i$  as function of the flow size distribution and the set of thresholds  $\alpha_i$ . A summary of all quantities that have been defined is provided in Table 2.1. These parameters are sufficient to express the average sojourn time on a single link, so to characterize the performance of PIAS. The average sojourn time on queue  $i$  which comprises the average waiting and serving time is given by the well-known formula for M/M/1 queues:

$$\mathbb{E}[T_i] = \frac{1}{\mu_i - \lambda_i} \quad (2.4)$$

The M/M/1 model holds for every subsequent queue. This follows from Burke's theorem [6], that states that the outgoing process of an M/M/1 queue is a Poisson process. Thus, cascaded queues still observe a Markovian arrival process.

Variable	Description
$Q_i$	Priority queue $i$
$N$	Number of priorities $i$
$X$	Flow size
$F(x)$	Flow size c.d.f.
$f(x)$	Flow size p.d.f.
$\lambda_i$	Packet arrival rate at PQ $i$
$\mu_i$	Drain rate of PQ $i$
$\rho_i$	Average load on $Q_i$
$L_i$	Customer size at PQ $i$
$T_i$	Waiting time at PQ $i$
$\alpha_i$	Demotion threshold from $Q_{i-1}$ to $Q_i$
$p_i$	Probability that a flow leaves after $Q_i$

Table 2.1: Variables of the model.

Finally, the total average sojourn time in the tandem of  $N$  queues is just a linear combination of  $\mathbb{E}[T_i]$ , where the coefficient that weight the individual sojourn times at any priority queue are the probabilities that a flow shall traverse the same PQ.

$$\mathcal{T} = \sum_{i=1}^N \theta_i \sum_{j=i}^N \mathbb{E}[T_j] \quad (2.5)$$

Thus, given the flow size distribution it is possible to easily compute the system performance according to the model yet presented.

### Optimal thresholds

Plainly, it follows that the *optimal* thresholds  $\alpha_i$  can be derived solving the non-linear minimization of  $\mathcal{T}$ . Notice that they are easily obtainable once all  $\theta_i$  are known as

$$\alpha_i = F^{-1}\left(\sum_{j=1}^i \theta_j\right)$$

Hence, it is possible to solve in  $\theta_i$ .

$$\begin{aligned} \min_{\{\theta_i\}} \quad & \mathcal{T} = \sum_{i=1}^N \theta_i \sum_{j=i}^N T_j \\ \text{subject to} \quad & \theta_i \geq 0 \quad \forall i \in [1, N] \\ & \sum_{i=1}^N \theta_i = 1 \end{aligned} \quad (2.6)$$

### 2.3.2 Greedy threshold assignment

In order to verify the actual gain obtained with optimal demotion levels, it is practical to compare it against simpler greedy strategies for thresholds assignment. Two intuitive

approach are considered for fast threshold computation: Equal Split (ES) and Load Split (LS).

### Equal Split (ES-N)

The *Equal-Split* method slices the flow size distribution uniformly in  $N$  percentiles. The resulting thresholds are the corresponding quantiles:

$$\alpha_i = F^{-1}\left(\frac{i}{N}\right), \quad i = 1, \dots, N-1$$

It is easy to observe that this criterion may be largely sub-optimal depending on the shape of the flow size distribution. In the case of very sharpened distributions short flows are demoted too early to lower priorities, while for heavy-tailed distributions, where high percentiles correspond to very long flows, latency sensitive flows remains mixed for long with elephant throughput-oriented streams.

### Load Split

*Load Split* is a technique where the thresholds are chosen in a way to control the amount of traffic fed in every priority queue. The problem is easy to understand in the simple case of  $N=2$  priority queues and a single threshold  $\alpha$ . Let  $G(y)$  be the amount of traffic generated by flows with size less than  $y$ :

$$G(y) = \int_0^y x f(x) dx$$

The traffic on high priority queue  $Q_0$  is derived as the sum of bytes generated by flows whose size is smaller than the threshold  $\alpha$ , and the bytes transmitted by flows larger than  $\alpha$ :

$$\mathbb{E}[L_0] = \int_0^\alpha x f(x) dx + \int_\alpha^\infty \alpha f(x) dx = G(\alpha) + \alpha \bar{F}(\alpha)$$

while the traffic on low priority queue is

$$\mathbb{E}[L_1] = \int_\alpha^\infty (x - \alpha) f(x) dx = \mathbb{E}[X] - \mathbb{E}[L_0]$$

Then, the traffic on the high priority queue can be controlled solving the following *load-balance* equation:

$$\mathbb{E}[L_0] = \gamma \mathbb{E}[X] \tag{2.7}$$

A proper choice of  $\gamma$  apportions to the high priority queue a fraction of the average total traffic  $\mathbb{E}[X] = \mathbb{E}[L_0] + \mathbb{E}[L_1]$ . For example, if a perfect load balancing between the two queues is desired then it is set  $\gamma = 1/2$ . Extending these equations to the general case of any number of priority queues is rather simple. The traffic on the  $i$ -th queue has the same expression of Eq(2.2). It can be rewritten shortly with the notation just introduced:

$$\mathbb{E}[L_i] = G(\alpha_i) - G(\alpha_{i-1}) + \alpha_i \bar{F}(\alpha_i) - \alpha_{i-1} \bar{F}(\alpha_{i-1})$$



At this point, similarly to the simple case of two queues, it is enough to solve the following *load-balance* equations iteratively for all  $i = 1, \dots, N - 1$ :

$$\begin{cases} \mathbb{E}[L_i] = \gamma_i \mathbb{E}[X], & i = 1, \dots, N - 1 \\ \sum_i \gamma_i = 1, & \gamma_i \in [0, 1] \end{cases} \quad (2.8)$$

In fact, also in this case it holds  $\sum_{i=1}^N \mathbb{E}[L_i] = \mathbb{E}[X]$ .

**Proof.** In the expansion of  $\sum_{i=1}^N \mathbb{E}[L_i]$  at each step new terms simplify old terms. Since  $\alpha_0 = 0$  and  $\bar{F}(\alpha_N) = 0$ , it gives  $\int_0^\infty x f(x) dx = \mathbb{E}[X]$ .

$$\begin{aligned} \sum_{i=1}^N \mathbb{E}[L_i] &= \dots + \cancel{G(\alpha_i)} - G(\alpha_{i-1}) + \cancel{\alpha_i \bar{F}(\alpha_i)} - \alpha_{i-1} \bar{F}(\alpha_{i-1}) + G(\alpha_{i+1}) - \cancel{G(\alpha_i)} \\ &\quad + \alpha_{i+1} \bar{F}(\alpha_{i+1}) - \cancel{\alpha_i \bar{F}(\alpha_i)} + \dots = G(\alpha_N) - G(\alpha_0) + \cancel{\alpha_N \bar{F}(\alpha_N)} \xrightarrow{0} - \cancel{\alpha_0 \bar{F}(\alpha_0)} \xrightarrow{0} \\ &= \mathbb{E}[X]. \end{aligned}$$

■

## Chapter 3

# The spatial-diversity framework

### 3.1 Improving scheduling with spatial-diversity

In the previous chapter have been addressed the scheduling of jobs of variable size in relation to their probability distribution. Then, there were presented two attempts to approximate the optimal LAS and SRPT schedulers by leveraging multiple priority queues at network interfaces. However, one limitation of this approach remains the scarce number of such priority queues available in commodity switches. In particular, it was mentioned that a large number of priority levels is demanded to better approximate the reference scheduling disciplines, in which  $N \rightarrow \infty$  (Sec. 2.1.2). Unfortunately, devices in modern DCNs are usually equipped with no more than 8 priority queues per port, whose majority are reserved for other purposes, like isolating different types of traffic. Indeed, several transport protocols may coexist in the same network, without necessarily being designed to be fair with respect to each other. Example of such transports are RDMA [9], DCTCP [2], standard TCP and UDP. For this reason, it is realistic to assume having at most  $N = 2$  priority queues in practical cases.

Upon this understanding, the main proposal of this work is to evaluate the possibility of exploiting the high degree of path diversity typically offered by DCN topologies, in order to improve the effectiveness of flow scheduling. Large DCN topologies usually are multilayer recursive Clos networks which offer a variety of equal-cost paths between racks. Precisely, in the simplest 2-layer Fat-Tree there are a number of paths proportional to the number of spines  $K$  (Sec. 1.1). The key observation is that much like priority queues, different paths yet are a way for augmenting the granularity of prioritization and for separating flows with different QoS requirements. The novel paradigm being investigated aims to exploit *spatial-diversity* to derive extra priorities and overcome the limitations in the maximum number of available PQs imposed by a single switch. Indeed, the mechanism of priority demotion, as proposed in PIAS, only shifts long flows across PQs of a single interface. Instead, we argue that the same demotion could be potentially applied at DC-level. Consider a Leaf-Spine topology and focus on a single ToR switch (Fig.3.1). The interfaces from such a ToR towards all spines can be seen jointly as a unique big interface with  $K$  times more priority queues than a port alone. The basic MLFQ (Fig.3.1a) focuses on the links individually and thanks to demotion moves flows, during their lifetime,

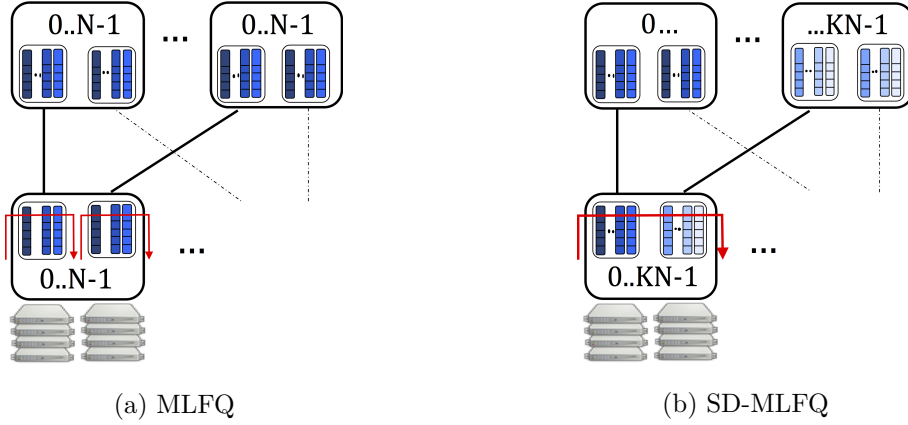


Figure 3.1: Demotion extension with spatial diversity. Red arrows indicate demotion trajectory of the longest flow. On each switch is reported the priorities it handles. Lower indexes (and darker colors) correspond to higher priorities.

across priority queues of single interfaces. Therefore, every interface handle the same priorities and flows are load balanced on different links independently from the prioritization mechanism. A standard technique for load balancing at flow-level is ECMP, which derives the next hop from the transport-layer tuple  $\{\text{IP ADDRESSES}, \text{PORTS}, \text{PROTOCOL ID}\}$ . Instead, our solution in Fig. 3.1b — which we call Spatially-Diverse MLFQ (SD-MLFQ) — takes advantage of spatial diversity to extend the number of demotion levels beyond the limitation imposed by the PQs on a single interface. Interfaces from any ToR to the connected spines are virtually aggregated to offer a wider range of demotion levels. One new aspect of this novel approach is that a demotion could imply shifting a flow from one path to another of equal-cost. In that light, the routing — thus load balancing — over the switching fabric is not blind to the prioritization machinery, but does depend on it. In general, a flow is moved both across queues and spines, effectively allowing a global resource exploitation for the demotion scheme. In order to have a clear and direct notation, the spines are labeled with the priorities handled by their interfaces. Notably, since the demotions that imply a spatial re-route take place at ToRs, all the interfaces on the same spine handle the same priorities. Hence, it is enough a single labeling per spine, valid for all its interfaces. A clear benefit of spatial diversity is that finer granularity in priorities is achieved even with few queues per port, at the price of a very limited implementation complexity. Also, elephant flows are better segregated from mice flows, as after a while they are physically moved to different paths inside the switching fabric. Despite its simplicity, relevant works exploring this solution in the field of data center networks seem to be lacking.

### 3.1.1 Abstraction as queuing system

To tackle the spatial diversity framework from a conceptual perspective, three queuing systems are compared. They are shown in Fig. 3.2. This is an abstraction of the real data center topology. At its heart, a Leaf-and-Spine network with  $K$  spines can be described as

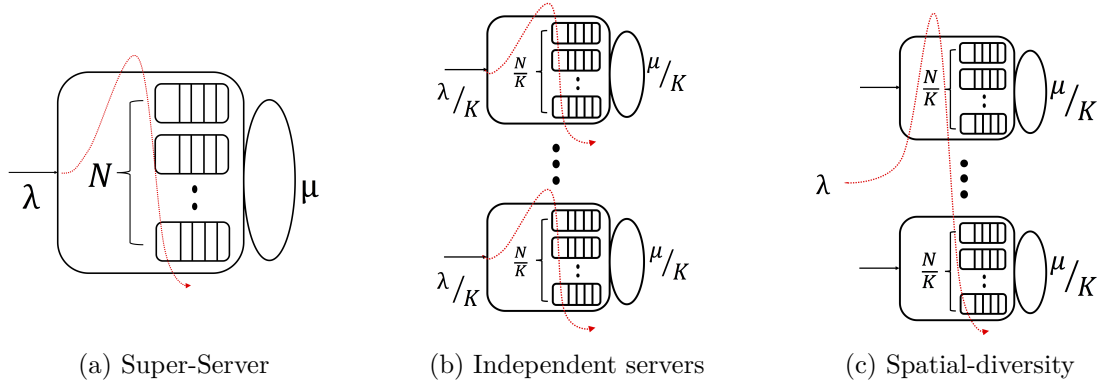


Figure 3.2: Three queuing system comparison.

a queuing system of  $K$  parallel G/G/1 servers. Actually, these servers can be equivalently mapped to the *up-send* interfaces that connects a ToR switch to all the spines, as well as to the *down-send* interfaces from different spines to a single ToR. This mapping with the leaf-and-spine topology is shown in Figure 3.3, where the up-send interfaces are the ones with red texture, whereas the spine down-send interfaces are the ones with light green texture. Ingress and egress ports that connect end hosts to the datacenter network are ignored at this stage. This is in contrast with the data center abstraction provided by pFabric as a giant switch (Figure 2.1), where the ingress and egress queues represented the bottleneck where to deploy scheduling strategies, whereas the switching fabric was assumed to be an ideal non-blocking interconnection. Somehow, differently from state-of-the-art solutions that focus only on the bottleneck links assuming that the network is able to sustain maximum throughput with negligible delays, the spatial-diversity approach shifts the attention to the queues inside the fabric itself and rethink the scheduling with a global DC view. Therefore, first it will be assessed the impact of spatial diversity with this simplified model which disregards the hosts and the ingress/egress interfaces, then we will proceed with the implementation on a simulated DCN afterwards.

From now on the interfaces considered in the abstract queuing model of Figure 3.2 will

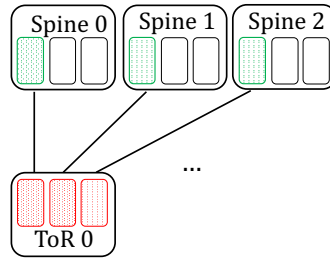


Figure 3.3: Mapping of the model abstraction with DCN topology

be interchangeably referred to as servers or spines (with reference to the mapping with down-send spine interfaces).

The three systems are different alternatives to handle the same total arrival rate  $\lambda$

with the same total processing capability  $\mu$  and the same number of priority demotion levels  $N$ . All systems use the mechanism of priority demotion first introduced by PIAS (§2.2.2). Instead, they differ in how priority levels are partitioned on a number  $K$  of parallel servers and, more importantly, in the way they use available servers for flow demotion. Specifically, the longest flow that experience all possible demotions follows different trajectory across priority queues and servers. Flow trajectories are represented with dashed red arrows in Figure 3.2

The first case (Fig. 3.2a) is to have a single high-capacity server that handles all priorities. This of course is the absolute best case where resources are fully concentrated, that provides the smallest delay but does not scale to the dimensions of real systems. It would be equivalent to realize an entire data center interconnection network with a single device of astonishing bandwidth. The second case (Fig. 3.2b) is the legacy way to handle priority demotion, where all servers are treated independently in parallel. Flows are evenly load balanced on the available servers and moved across priorities of the same server during their lifetime. In this case the number of demotion levels are limited to  $\frac{N}{K}$ . Finally, the third case (Fig. 3.2c) is the novel object we want to investigate, where all the flows are initially sent to the same server, then demoted on the  $N$  globally available priority queues. In this case subsequent servers are configured to handle lower and lower priorities, thus part of the flows are re-routed as a consequence of the spatial-diversity. For this reason, in the following of this work also the servers and the links will be labeled as "high priority" or "low priority" for brevity, meaning that they handle high priority traffic or low priority traffic, respectively.

It was already defined a rigorous mathematical formulation for representing the demotion across priority queues in a single servers as a tandem of M/M/1 queues. (§. 2.3.1). The system with spatial diversity introduces a new set of thresholds, which mark the amount of service after which a flow is rerouted to another server. The next section presents a complete formulation for modeling spatial-diversity, starting from the model already defined.

## 3.2 Mathematical formulation for spatial-diversity

The introduction of spatial diversity adds a level of complexity to the system. Let's formalize the system setup. There are  $K$  parallel servers with  $N$  priority queues each one. Thus, there are a total of  $K \times N$  priority queues. All of them are used to add new levels for demotion, therefore flows can assume priority  $p \in [1, KN]$  with no duplicate priorities and there are  $K \times N - 1$  demotion thresholds. Let's use the variable  $j \in [1, K]$  to index the servers and the variable  $i \in [1, N]$  to index the priority queues inside each server  $s_j$ . Starting from the highest level of priority  $p=1$  and following descending order of priority up to  $p = KN$ , the priority levels are assigned to servers from the lower index to the higher index. Thus, server  $s_1$  handle priorities  $p = \{1, \dots, N\}$ , server  $s_2$  handle  $p = \{N+1, \dots, 2N\}$ , and so on. As in the MLFQ system, whenever a flow at priority  $p$  has received service equal to the next threshold, it is downgraded at priority  $p+1$ . However, differently from MLFQ, it may happen that priority  $p$  is assigned to server  $s_j$ , while  $p+1$  is handled by server  $s_{j+1}$ . In such a case, the flow must be rerouted to a different server.

Thus, two related problems need to be solved: finding the set of thresholds that triggers a new flow route from server  $s_j$  to server  $s_{j+1}$ , and finding for each server  $s_j$  the set of thresholds that mark the demotion from PQ  $i$  to PQ  $i + 1$ , when both PQs are in server  $s$ . Let's denote them as *load-balance thresholds* and *sub-thresholds*, respectively.

#### Terminology

In a system with priority demotion and with spatial diversity, a flow is moved to a different priority queue (PQ) when the service it has obtained is larger than some threshold. We say that the threshold *push* a flow to a new PQ. We recognize two set of thresholds:

- ***Load balance thresholds.*** The set of demotion thresholds that push a flow in a priority queue on a server different from the server currently handling the flow itself, implying a flow rerouting.
- ***Sub-thresholds.*** The set of thresholds that push a flow in a different priority queue but in the same server as the one currently handling the flow itself.

Correspondingly, we will refer to as:

- ***Inter-server[spine] demotion.*** A priority demotion involving a load balance threshold.
- ***Intra-server[spine] demotion.*** A priority demotion involving a sub-threshold.

The name *load-balance* derive from the fact that the values of these thresholds have strong implications on how the load is distributed across servers. If the thresholds are too small, flows are early rerouted on low priority paths and the capacity of high priority links is essentially wasted. For a real data center implementation this would mean a reduction in the maximum throughput sustained by the switching fabric. As a trivial example, consider a simple topology with two parallel servers each one of normalized capacity 1. In this topology there is only a single threshold that may trigger a flow reroute. With an ideal load balance that evenly distributes the traffic among the two servers, this topology offers a maximum normalized throughput of 2. However, an inappropriate setting of the aforementioned demotion threshold that immediately reroutes flows after negligible attained service would in practice reduce the total capacity of 50%. On the other end, still this threshold could be optimal if considering only micro flows, whereas the equal load balance which uniformly splits the traffic on available servers may correspond to a bad demotion threshold for the FCT minimization. In short, a careful tuning of load balance is a new trade-off that was nonexistent in the legacy MLFQ framework, where traffic was demoted to different priority queues but always in the same interface (i.e. link), with zero effects on load balancing. Indeed, a threshold setting that gives an unbalanced load allocation in the available priority queues affects only the delays but not the maximum throughput that the network can sustain.

As concerns the sub-thresholds, they push flows to the same kind of intra-server demotion already studied for MLFQ. However, MLFQ servers work independently and in parallel (Fig. 3.1a) and they are all fed with the same job size distribution. Instead, SD-MLFQ servers observe a different version of the initial workload. Indeed, subsequent servers receive only flows larger than the precedent demotion thresholds, thus they handle truncations of the original flow size distribution above increasing percentiles. As a consequence, while MLFQ thresholds are computed once for all servers, SD-MLFQ sub-thresholds should be optimized depending on the server they belong to. A clarifying overview is provided in Fig. 3.4 for the case of  $K=3$  servers  $s_1, s_2, s_3$ . Figure 3.4a is an example job size distribution, from which new flows are randomly generated. In a data center network, this would be the DC-wide workload. On the same axes are drawn the two example split thresholds, denoting the amount of service in kilobytes after which a flow is rerouted from  $s_1$  to  $s_2$  and then from  $s_2$  to  $s_3$ . Since all flows enter the system

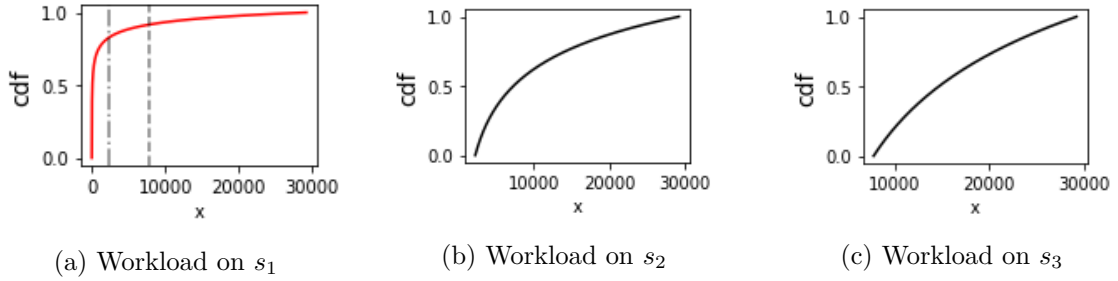


Figure 3.4: Job size distributions observed by servers  $s_1, s_2, s_3$  in decreasing order of priority. Optimal load balance split for  $\lambda = 0.9$ .

through  $s_1$ , this distribution is also the one observed by the highest priority server  $s_1$ . Then, subsequent servers of lower priorities observe left-truncated versions of the initial workload on  $s_1$ .

An intuitive solution to address both load balance thresholds and sub-thresholds at once is to write an extension to the PIAS model that embeds spatial-diversity. This formulation, explained in details below, would in principle guarantee optimal performances. Indeed, it jointly captures in a unique model all the dynamics of the system.

### Optimal solution

The optimization problem Eq.(2.6) can be extended to the case of spatial diversity. Assume the capacity of a single server to be  $\mu/K$ , so that the total capacity is always  $\mu$ . All the servers work independently with same rate  $\mu/K$ . Still there is a SP scheduler orchestrating priority queues in every server, but there is no coordination among different servers, meaning that there is not a global SP scheduler to discipline the transmission among distinct servers. Indeed, such a tight control among physically different devices would be almost impossible to realize in practice, at least packet by packet.

Let the notation  $\mu_i^j$  indicate the drain rate of the  $i$ -th queue in the  $j$ -th server  $s_j$  ( $1 \leq i \leq N, 1 \leq j \leq K$ ), and equivalently use  $\lambda_i^j$  for the arrival intensities and in general the superscript  $j$  for any quantity already defined for the queues in the basic

model in Sec. 2.3.1 but applied to queues inside server  $s_j$ . A summary of all the quantities involved in the queuing model is found in Table 2.1. It follows:

$$\begin{aligned}\lambda_i^j &= \lambda \mathbb{E}[L_i^j] \\ \mu_i^j &= \frac{\mu}{K} \prod_{l=1}^{i-1} (1 - \rho_l^j) \\ \theta_i^j &= F(\alpha_i^j) - F(\alpha_{i-1}^j) \\ \mathbb{E}[T_i^j] &= \frac{1}{\mu_i^j - \lambda_i^j}\end{aligned}$$

The formulation is essentially the same but with an additional dimension that takes into account the existence of multiple servers. Every server is modeled as a tandem of  $N$  queues, exactly as in PIAS. Then, there are  $K$  of them in parallel. Also, we wrote an additional constraint Eq.(3.1d) that prevents overloading any spine (TODO: additional constraint is it really needed??? in case of overloading, sojourn times goes to infinity, so it is unlikely it is chosen as solution....).

$$\min_{\{\theta_i^j\}} \quad \mathcal{T} = \sum_{j=1}^K \sum_{i=1}^N \theta_i^j \sum_{j=i}^N T_j^s \quad (3.1a)$$

$$\text{subject to} \quad \theta_i^s \geq 0 \quad (3.1b)$$

$$\sum_{i=1}^N \theta_i^j = 1 \quad \forall j \in [1, K] \quad (3.1c)$$

$$\sum_{i=1}^N \lambda_i^j < \mu/K \quad \forall j \in [1, K] \quad (3.1d)$$

This approach would provide optimal load balancing and at the same time would choose a set of sub-thresholds targeted on the optimal load balance thresholds. Notably, the optimal solution would never include load balance thresholds that lead to an unbalanced traffic distribution among the servers, because that would have an overkilling effect on the delays, which instead the problem tries to minimize. Despite being a clean analytical formulation, its complexity seems prohibitive. The basic model without spatial diversity yet was non-convex and presented products and ratios of variables (Sec. 2.3.1). Nonetheless, it is still tractable since the number of variables is typically bounded to the (low) number of PQs available in commodity switches. Instead, in the complete model the number of variables scales with the product  $K \times N$ , where  $K$  is usually big for large-scale data centers. In the next chapter we will provide in more details the CPU-time spent by two well-known meta-heuristics solvers, specifically PSO [10, 14] and Basin-Hoppin [15] before to converge. Anyway, if the time scale at which a solution could be found is excessively large, the system would be unable to react promptly to a sudden change of the flow size distribution. As a consequence, the system would operate for long using thresholds mismatched with respect to the flow distribution. Therefore, even finding an optimal solution once does not signify that this approach is feasible for a real datacenter scenario, where likely the solution must be computed repeatedly as statistics change.



To the purpose of investigating the spatial diversity framework, it has been preferred to handle the two problems individually. The approach that has been carried out is decoupled in two sequential steps.

1. **Optimize load-balance thresholds.** First it is solved an optimization problem with the goal of finding the best load partitioning among  $K$  servers. Servers are always assumed to have a single priority ( $N=1$ ), because for the time being sub-thresholds are ignored. At the end of this phase a set of  $K-1$  load balance thresholds is delivered.
2. **Greedy subthresholds.** The load-balance thresholds are provided as input of the second phase. They split the support of the flow size in  $K$  disjoint intervals covering the whole support. The  $i$ -th interval contains the sizes of those flows that end their service in server of priority  $i$ . On each interval is computed a set of sub-thresholds with a greedy algorithm like ES-N or LS-N, but applied to the truncated distributions (Sec.2.3.2).

### 3.2.1 Decouple load balance from sub-thresholds

#### Optimize load balance thresholds

For the first phase turns out again to be useful the stochastic queuing model of MLFQ as provided in PIAS. (Sec.2.3.1-Fig.2.2). Each server is equipped with a single priority queue per port, therefore the optimal load balance problem can be abstracted with exactly the same model, where each queue in the tandem maps a server. After all, in this case priority queues are physically distributed to different servers, instead of being part of the same interface. They are independent on each other and the strict priority scheduler is no more involved, as they work in parallel without coordination. Practically, the only modification is to the queues capacities  $\mu_i$ . Remember that in the aforementioned MLFQ model the strict priority scheduler was described by attenuating the PQ serving rates  $\mu_i = \mu \prod_{i=1}^{i-1} (1 - \rho_i)$  for increasing  $i$ . Because of servers work in parallel without scheduling, this expression is not needed anymore and the draining rates just coincide to the same value:

$$\mu_i = \mu / K, \quad i = 1, \dots, K$$

#### Choose greedy sub-thresholds

Lastly, it is addressed the problem of finding the sub-thresholds. Starting from the optimal load balance thresholds that have been found in the previous step, each server is then treated individually. Remember that whatever policy is adopted for sub-thresholds computation, it has to be applied to all servers individually, since they observe different flow size distributions. All the threshold computation algorithms defined so far (§2.3.1, §2.3.2) require the knowledge of the flow size distribution. It is pretty straightforward to obtain the p.d.f.  $f(x)$  and the corresponding c.d.f.  $F(x)$  as distribution conditioned to the truncated supports. There are  $K$  servers and  $K - 1$  load balance thresholds. For the sake of conciseness, denote the load balance thresholds  $\alpha_N^j$  which triggers a flow reroute

from server  $j$  to server  $j + 1$  as  $\Omega_j$ . On server  $j$  the initial distribution is normalized on a new support  $[\Omega_{j-1}, \infty)$ . Thus, from probability theory:

$$\begin{aligned} F(x|x > \Omega_{j-1}) &= \frac{F(x) - F(\Omega_{j-1})}{1 - F(\Omega_{j-1})} \\ f(x|x > \Omega_{j-1}) &= \frac{f(x)}{1 - F(\Omega_{j-1})} \end{aligned} \quad (3.2)$$

Once these distributions are known, nothing is missing to apply the greedy sub-threshold assignment algorithms depicted in Sec. 2.3.2. For completeness, they are briefly rewritten under this framework. Define in short:

$$\begin{aligned} f_T(x) &= f(x|x > \Omega_{j-1}) \\ \bar{F}_T(x) &= \bar{F}(x|x > \Omega_{j-1}). \end{aligned}$$

Enumerate the available server  $s \in \{1, \dots, K\}$ . The goal is to find on all servers a set of sub-thresholds  $\alpha_i^s$  to delimit demotion bounds across priority queues  $i \in \{1, \dots, N\}$ .

*Equal-Split-N* (ES-N) is substantially the same:

$$\alpha_i^j = F_T^{-1}\left(\frac{i}{N}\right) \quad (3.3)$$

*Load-Split-N* (LS-N) is also very similar. It required the solution of the set of load balance equations 2.8. The expression of the average traffic  $\mathbb{E}[L_i^j]$  on priority  $i$  of server  $s_j$  is unchanged. However, additional care must be paid to their sum  $\sum_{i=1}^N \mathbb{E}[L_i^j]$ . From proof 2.3.2, this sum was equal to:

$$\sum_{i=1}^N \mathbb{E}[L_i^j] = \int_{\Omega_{j-1}}^{\Omega_j} x f_T(x) dx + \Omega_j \bar{F}_T(\Omega_j) - \Omega_{j-1} \bar{F}_T(\Omega_{j-1})$$

and the terms outside the integral always amounted to zero, so that the sum gave  $\int_{\Omega_{j-1}}^{\Omega_j} x f_T(x) dx = \mathbb{E}[X]$ . Instead, in this case they are generally different zero because the lower threshold  $\alpha_0^j = \Omega_{j-1} \neq 0$  and the survival function above the upper threshold  $\bar{F}_T(\alpha_N^j) = \bar{F}_T(\Omega_j) \neq 0$ . Everything else is the same.

As already mentioned, it is very likely that this approach is sub-optimal. However, it represents an handful way of computing all the thresholds in order to establish some understanding about the integration of spatial diversity with an MLFQ system.



## Chapter 4

# Numerical analysis of spatial diversity

The previous chapter introduced the idea of exploiting spatial-diversity to provide more priority levels for flow classification. It was discussed why in principle such technique could yield flow completion time gains, especially when few priority queues are available. It was explained that adopting a spatial-diversity demotion has strong implication on the traffic load balancing, since flow routing becomes priority-dependent. A few questions arises soon: how to choose the thresholds to distribute the load on the topology? What are the relationships between the priority granularity in single interfaces and spatial diversity? How to scale spatial diversity with the topology size? The goal of this chapter is to validate our intuitive hints with numerical results and to shed the light on the benefits and the restraints of the proposed algorithm. In particular, it will go through an exhaustive analysis of the system by delivering plenty of numerical experiments obtained with a custom flow level simulator implemented in Python. This simulator does not capture any of the complex dynamics inherent to a real packet network, it does not have any protocol stack implemented neither at traffic sources nor in the switching modules. Rather, it is a job-oriented queuing simulator that disregards packet level events but only runs flow arrivals and serve them in generic queues. Its purpose is to provide a clean baseline numerical analysis not plagued by possible side effects due to network misconfiguration. Next sections will be an in-depth analysis of dynamics of the spatial diversity, when varying the dimensionality of the system both in the number of servers and in the number of priorities.

## 4.1 Model implementation

### 4.1.1 Workloads

The performances of the three systems are compared using two empirical flow size distributions that have been derived from production data centers (Figure 4.1). Flow size distribution is shortly termed *workload*. The first workload has been estimated instrumenting thousands of servers in a datacenter hosting a Web search [2] application, while

the other refers to data mining tasks [8]. As expected, these distributions have a mix

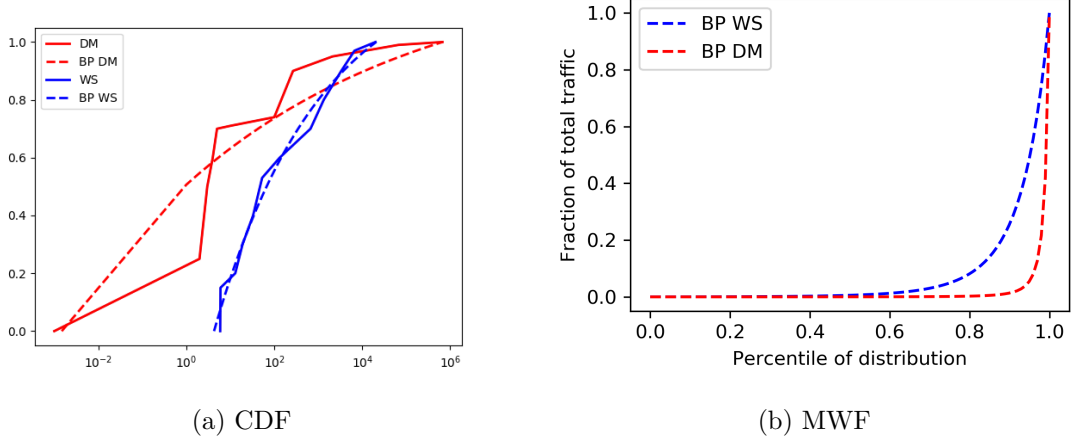


Figure 4.1: Workload properties

of short and long flows and both present the high-variability property typical of data center traffic (§1.2.2). Figure 4.1a shows with solid lines the cumulative density function of the two empirical workloads, along with two analytical bounded Pareto distributions (dashed lines), whose parameters have been fitted to the corresponding empirical points. The bounded Pareto distribution is a truncated version of the Pareto distribution over the finite support  $[u, t]$  and it is well-suited to model heavy-tail characteristics. It has three parameter: the lower extreme of its support  $u$ , the upper extreme  $t$  and the shape parameter  $\alpha$  that controls the weight of its tail. The analytical expression of its cdf  $F(x)$  in the interval  $[u, t]$  is:

$$F(x) = \frac{1 - \left(\frac{u}{x}\right)^\alpha}{1 - \left(\frac{u}{t}\right)^\alpha}, \quad 0 \leq \alpha \leq 2 \quad (4.1)$$

This distribution has been chosen to be used in the analysis due to some graceful properties. First of all, it is relatively easy to control its variability by a proper tuning of its parameter  $\alpha$ . Values of  $\alpha$  close to 2 accentuate the heavy-tail property, while smaller values of  $\alpha$  tend to regularize a bit its behavior. Second, being definite on a limited support it can be adapted to any minimum and maximum flow size in the datacenter. Third, the Pareto distribution is scale-invariant, meaning that normalized Pareto distributions remains Pareto. Nicely, this implies that the workloads observed by subsequent servers can be always modeled with the same probability distribution, only changing parameters. After all — as seen in chapter 3 — these workloads are conditional distribution obtained with simple normalizations. Last, its mean and its variance — which depend on  $\alpha$  — are finite, thus the problem of finding the shape parameter for any fixed first and second moment can be smoothly treated numerically. Specifically, for the bounded Pareto, the

mean and the variance have the following expression:

$$\mathbb{E}[X] = \frac{\alpha}{(1 - \alpha)(t^\alpha - u^\alpha)}(u^\alpha t - t^\alpha u)$$

$$\sigma_X^2 = \frac{\alpha}{(2 - \alpha)(t^\alpha - u^\alpha)}(u^\alpha t^2 - t^\alpha u^2)$$

The best fit to the empirical distributions has been obtained with a simple Maximum Likelihood Estimator (MLE). The resulting parameters are reported for both the workloads. Let  $X$  be the flow size random variable as usual, and write in short notation  $BP(u, t, \alpha)$  the bounded Pareto.

$$X_{WS} \sim BP(3, 29000, 0.125)$$

$$X_{DM} \sim BP(0.1, 100000, 0.26)$$

The measurement unit for the extremes of the support in this case is kilobytes. The fitting error is higher for the data mining workload than for the web search. For low percentiles this is difficult to avoid because a very crude sampling is provided. In fact, on a total of 11 empirical points, 4 of them are for values above the 90th percentile. Instead, for high percentiles a better fitting likely could be obtained by weighting more the tail of the distribution.

The rightmost plot (Fig. 4.1b) completes the picture by showing the Mass-Weighted Function  $M_w(x)$  [7]. This can be seen as the probability that a byte picked at random belongs to a flow below a given percentile and it is used to characterize the variability of a distribution. Its name comes from its definition, where job sizes are weighted by their probability mass:

$$M_w(x) = \frac{\int_0^x x f(x) dx}{\mathbb{E}[X]}$$

It holds:

$$\int_0^x x f(x) dx \leq \mathbb{E}[X]$$

In other words, it is just the average normalized traffic injected by flows shorter than  $x$ . The figure has on the abscissa the percentile rather than the corresponding job size, to allow the comparison between workloads with different supports on the same axis. If  $y$  is a given percentile, it is evaluated  $M_w(F^{-1}(y))$ .

In summary, both distributions exhibit high variability. In the web search case the largest 4% of flows carry half of the total traffic, the data mining is even more skewed: 70% of the flows are less than 8 packets only, but almost the entire load is sustained by a ridiculous percentage of flows of about 100MB of size. This suggests that the more challenging distribution to schedule is the web search, consequently it is the one that will deserve most of the attention. In fact, recall that an ideal flow-agnostic LAS scheduler guarantees lower and lower delays as the variability of the distribution increases, both on average and at high percentiles (§2.1.2). The theory is confirmed pretty straightforwardly by the simulation results presented next. Moreover, the web search distribution is also a lot easier

to simulate, since the very long tails of the data mining workload require protracted time-consuming simulations before being precisely reproduced. Long flows occur sporadically, however they give the main contribution to generate a desired load on the system.

#### 4.1.2 Optimal traffic load balancing

In previous discussions, it was already realized that the spatial diversity framework introduces strong implications on how the load is distributed on the switching fabric. We decided to solve an optimization problem with the goal of finding the optimal load balancing and then to treat sub-thresholds a posteriori (§3.2.1). Thus, we considered the abstraction of spatial diversity as a queuing system (§3.1.1) setting  $N=1$  priority queue per server. With this setup all flow demotions correspond to shifting a flow from one server to the other. This way, the original problem of jointly optimizing inter and intra server demotion at once, has been simplified to finding the load balance that minimize the average flow completion time.

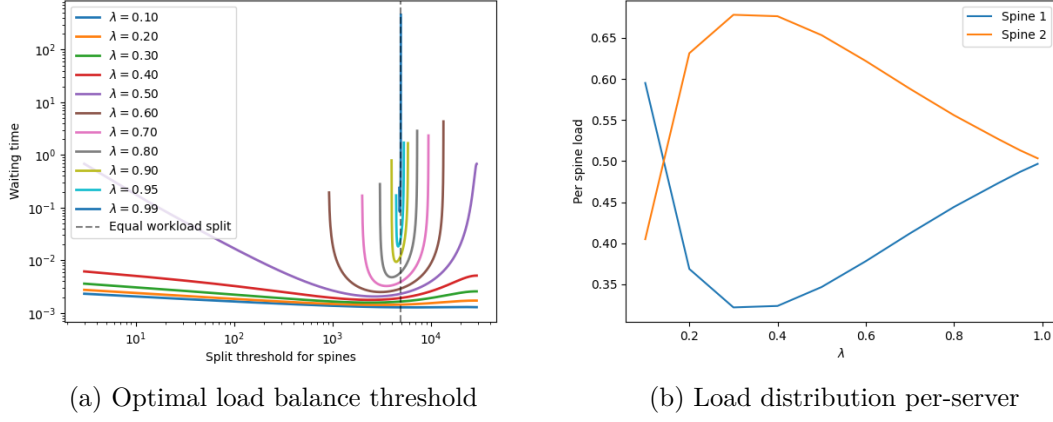
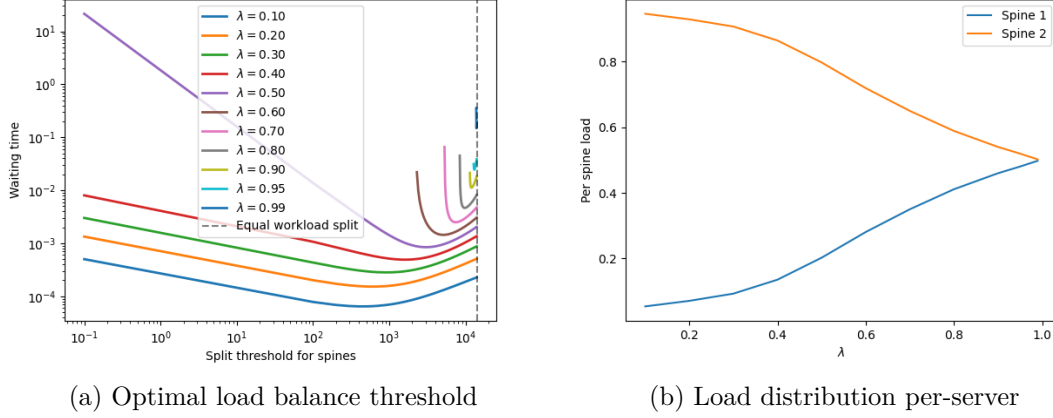
In this section we first start by analyze the simpler example of spatial diversity, where only two M/M/1 servers are deployed in parallel, each of them with only a single priority queue. In this case there is globally only one threshold, therefore it's possible to plot the shape of the cost function and to study its properties. It is worth remarking again that this threshold is a load balance threshold, therefore the following analysis will refer to the load balancing minimization problem, where there isn't any strict priority scheduler and all servers work in parallel with full capacity  $\mu/K$ . In other words, there is not any throttling that would assign only the residual capacity to the low priority server (already discussed in §3.2.1). We start initially by solving the formulation for M/M/1 queues, then we will look for the solution of M/G/1 as well. For the sake of simplicity, we have implemented a total service rate  $\mu = \mathbb{E}[X]$ . In this way the average load fed in the system

$$\rho = \frac{\lambda}{\mu} \mathbb{E}[X]$$

is given only by the flow arrival intensity  $\lambda$ .

#### Load balance on 2 parallel servers

In the basic case of  $K=2$  parallel servers, it is possible to plot the average sojourn time when varying the load balance threshold. Figures 4.2a and 4.3a show how the cost functions look like, for the web search and the data mining workloads, respectively. All the axes are in log-scale and each curve represents a different normalized traffic  $\lambda$ . The vertical lines correspond to the split that gives perfect load balance, apportioning half of the traffic on the high priority server and half on the low priority one. This split may be interchangeably referred to as *perfect split* or *proportionate split* in the following. Figures 4.2b-4.3b show in parallel the normalized traffic distribution on the two servers, corresponding to the optimal threshold. One phenomena is visible for both workloads, confirming previous intuitions. The optimal load balance threshold does not coincide, broadly speaking, with the proportionate split threshold. Depending on the traffic level at which the system is operated and the workload, the optimal threshold triggers an earlier or


 Figure 4.2: Web search workload. Simple case of  $K=2$ .

 Figure 4.3: Data mining workload. Simple case of  $K=2$ .

later demotion with respect to the perfect split case. Equivalently, the jobs are distributed unfairly between the two servers (Figures 4.2b-4.3b). Imagine to connect the absolute minimums of the cost functions. For data mining, the imaginary line would be always in the leftmost side with respect to the proportionate split. That is, apart from loads close to saturation, the high priority server is kept as jobless as possible and the majority of work is sent to the low priority server. Remember that the problem formulation weighted the average sojourn times in the  $i$ -th priority queue with the percentage of flows with size between the thresholds  $\alpha_i$  and  $\alpha_{i+1}$ . The objective was:

$$\mathcal{T} = \sum_{i=1}^N \theta_i \sum_{j=i}^N T_j$$

Since the data mining workload is highly dominated by short flows, they receive more importance and longer flows are moved soon to another path. Plenty of short flows



carrying few bytes are kept on a separated link from medium sized flows and a couple of long flows. Thus, the result is not a surprise but it is coherent with the flow distribution. Slightly different trend is observed for web search. Moving from low to high values of  $\lambda$ , the optimal threshold is initially greater than the perfect split axis, then switches to its left and finally the two coincide. This is clear from the load distribution on the two servers. In general, for web search the load remains more balanced between the two servers in respect to data mining. The more unbalanced split occurs at  $\lambda=0.3$  where 70% of the traffic is handled by the low priority server. Instead, data mining has much more extreme load subdivision, especially for  $\lambda=0.1$ . This is a consequence of the high variability of the workload: there is six order of magnitude difference between the shortest and the longest flow and there is a pronounced heavy-tail. Hence in order to have significant changes on the weights  $\theta_i$  the threshold is moved significantly along the heavy-tail. In other words, the optimal threshold reroute few flows but lot of traffic on the low priority server. In fact, for  $\lambda=0.1$  the absolute value of the optimal threshold is much smaller than the proportionate split threshold, however only a small fraction of flows fits into this gap. Also, for similar reasoning, the two workloads have different sensitivities to the load balance threshold optimization. In particular, the web search achieves appreciable FCT gains starting from medium loads only ( $\lambda > 0.5$ ) with (TODO provide some numbers), whereas at low loads there is no practical difference with the proportionate split. On the contrary, the data mining distribution gives theoretically a (TODO provide some numbers) lower waiting time even for  $\lambda=0.1$ . Finally, for both cases the objective function becomes much more extremely curled and steep around the perfect split axis when the load approaches the saturation value 1. As already remarked, this is inevitable in order to fully exploit all the available capacity offered by the two parallel servers. At so high load any other traffic split would overload one of the two links and strongly deteriorate the average completion time of the flows going through it. Consistently, the cost function grows very rapidly in the neighborhood of the perfect split threshold value. Indeed, the M/M/1 response time  $T$  grows exponentially close to saturation. Its law (Figure 4.4) is:

$$T \propto \frac{1}{1 - \rho}$$

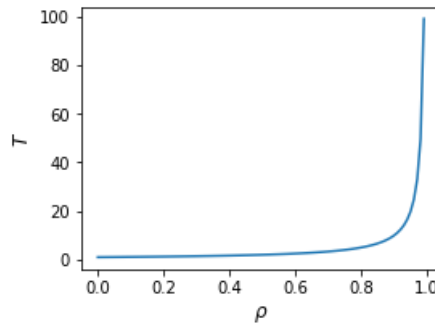


Figure 4.4: M/M/1 response time

### Validity of the model

Summarizing, the above discussion confirms that the web search traffic is harder than the data mining to cope with. This is the reason why in the following some results are shown exclusively for this workload. Next, it is consolidated the validity of the stochastic optimization model with numerical simulations, in particular it is shown the effective benefit of optimal load balancing. The underlying topology is again the simplest one, comprised of two parallel M/M/1 servers with no inner prioritization. The serving discipline is Processor Sharing (PS), implemented with a fluid model where parallel flows are served with equally subdivided bandwidth. The average normalized flow completion time (nFCT) has been considered as the primary evaluation metric.

#### Definition

**nFCT:** Given a fixed data center topology and pair of source-destination servers,  $\{s, d\}$ , define  $FCT_{opt}(x)$  as the FCT achieved by a flow of length  $x$  originated from  $s$  and directed to  $d$  in a completely empty DCN at load zero (excluding such a flow). Let  $FCT(x)$  be the FCT of a flow of length  $x$  in a DCN in the presence of other flows. Let  $X$  be the set of all possible flow lengths. Define:

$$nFCT = \sum_{x \in X} \frac{FCT(x)}{FCT_{opt}(x)}.$$

The normalized flow completion time is very similar in spirit to the average slowdown presented in Sec.2.1.2, but it is more handful as it is dimensionless. Its advantage is to put all flows on the same comparable scale, permitting a clean visual analysis of fairness with respect to flow size. By the way, this kind of evaluation is important to us, as spatial diversity mainly targets short flows.

Figures 4.5-4.6 report the nFCT gain obtained thanks to the sole load balance optimizer.

Two ( TODO three if we talk about M/G/1 ) different scenarios are compared. Both

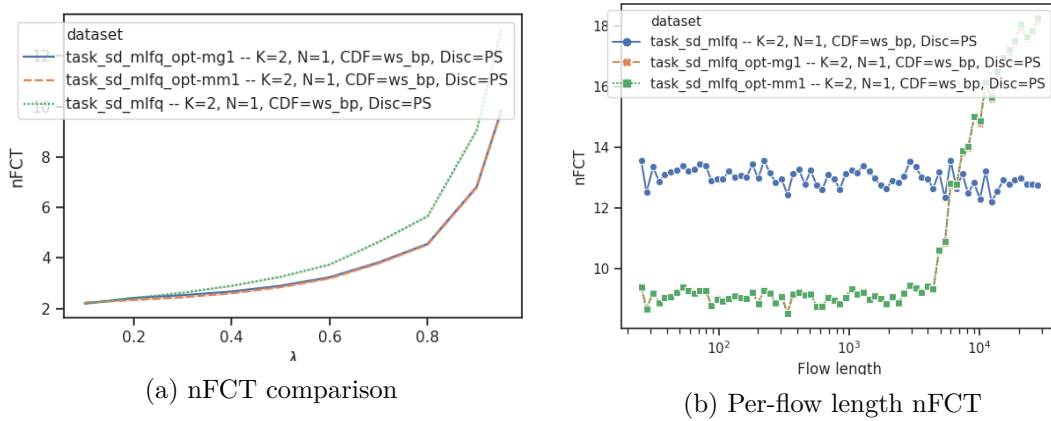


Figure 4.5: Web Search workload

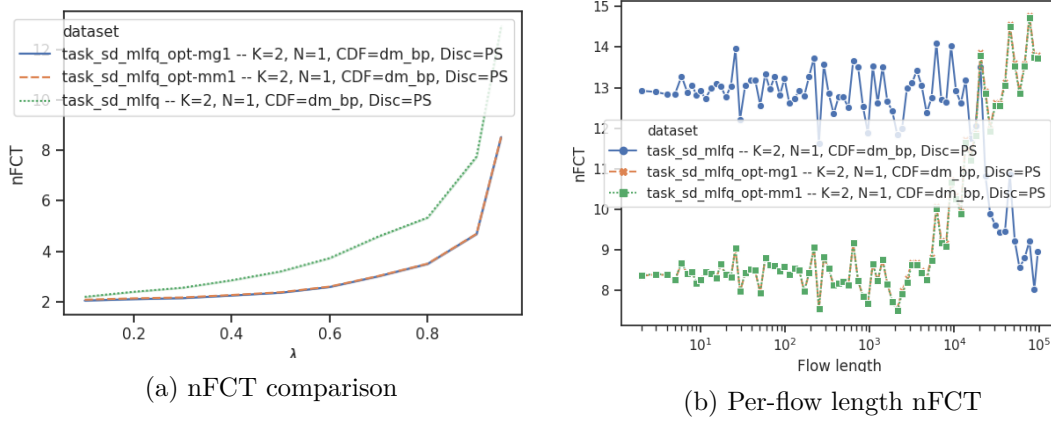


Figure 4.6: Data Mining workload

adopt spatially-diverse MLFQ, but in one case the load balance threshold is optimized, in the other not. The former case is labeled as **SD-MLFQ-OPT**, whereas the latter as **SD-MLFQ**. Call  $s_0$  and  $s_1$  the two servers (i.e links/queues) and  $\Omega$  the unique demotion threshold. Newly arrived flows always enter the system through  $s_0$  as their attained service equals to 0, then they are rerouted to the link  $s_1$  when  $s_0$  has transmitted  $\Omega$  of their bytes. As expected, the benefits on the average completion times are more pronounced for the data mining distribution, at all loads (Figures 4.5a-4.6a). Even so, appealing phenomena emerge when looking at the detailed breakdown of the response time versus flow size (Fig.4.5b-4.6b). First focus on the SD-MLFQ system without load balancing optimization and compare the blue curves. In web search, the normalized response time curve exhibits a constant plateau which testify a substantial invariance of the slowdown with respect to the flow length. Conversely, in data mining there is an abrupt transition to lower response times for flows longer than 20MB, which roughly corresponds to the demotion threshold adopted in this case. The difference is justified by the fact that the proportionate split threshold for data mining cuts the flow size CDF near the 99-th percentile, due to the oft-repeated heavy-tail characteristic of this workload. Therefore, only few long flows remain to share the processor of the low-priority server  $s_1$  — keep in mind that all servers have PS. Since the normalized FCT captures the slowdown in respect to the ideal case where the flow is serviced alone, it is explained such a behavior. Next, let's concentrate on the spatial diversity with optimal load balancing (green curves), which is the real goal to validate the model. The two workload perform similarly and have analogous trends. Specifically, the encouraging result is that shorter flows indeed take advantage of spatial diversity, as they are sent through a less loaded path. The opposite happens to long jobs, which undergone the dual effect.

### Analysis of the model scalability

So far the only scenario that has been investigated is the simple case of two parallel servers without sub-thresholds. This was useful to get confident with the relationships

between load balance and the spatial diversity demotion mechanism. However, in the more general setting there might be much more parallel servers and at least two priority levels per server. The total number of thresholds was reduced by decoupling the load distribution by the server-local prioritization. Thus, the total complexity only depends on the number of servers  $K$ . Unfortunately, it is still typically high. **TODO basin hoppin brief description (GERMAN) + parameter setting of the solver (GERMAN) + solution time plot description (ALE)**

## 4.2 The service discipline: FIFO or PS?

In the previous section it was laid down the basic framework of spatial diversity. It was remarked the fact that adopting spatial diversity translates into a priority-driven load balancing. At the same time were collected encouraging results about the effectiveness of the optimal load balancing. It was discussed the impact of two workloads commonly used as a benchmark in literature. Then, it was solved the optimal load balancing problem for up to  $K=9$  parallel servers and explained how to compute the sub-thresholds without trying to face the overwhelming complexity of the optimal approach. In all the simple experiments carried out, the servers were always configured with Processor Sharing service discipline and without priority queues.

This section aims to provide answers to principally two things. The first one is the behavior of the priority-driven load balancing when increasing the number of server  $K$  up to 9 parallel servers. Second, it is treated the integration of the priority-driven load balancing with the legacy MLFQ system. In practice, it is considered the general case where each server has many PQs on its interfaces used in strict priority. As it will be explained shortly, it turns out that for both attempts, the performances of the system are highly correlated with the adopted servicing policy. Considered cases are the FIFO (FCFS) servers and the usual PS servers.

### 4.2.1 Dimensioning spatial-diversity

First it is addressed the first question: what are the effects of augmenting the spatial diversity *rank*?

**SD-rank:** Given a data center topology with  $S$  spines — or the equivalent system of parallel M/M/1 servers where is applied priority-driven load balancing, that is each server handle a group of priority and a flow is routed depending on its assigned priority. Let's define the spatial-diversity rank (SD-rank) as the number  $K$  of servers which handle different priorities.

The spatial diversity ranks tells how many M/M/1 servers (or group of servers) are considered to be used to implement spatial diversity. For example, suppose there are 4 parallel servers available  $s_0, s_1, s_2, s_3$ . Consider these two possibilities.

1. **Rank 2.** The four servers are grouped in two pairs, say  $(s_0, s_1), (s_2, s_3)$ . Spatial diversity demotion is applied at pair level. This means there is only one load

balance threshold  $\Omega$ . When the longest flow enters the system, it is sent at random either on  $s_0$  or  $s_1$ . As soon as the flow has received service  $\Omega$  it is demoted, again randomly, either on  $s_3$  or  $s_4$ . In other words, this can be seen as a system of only two parallel servers with twice the capacity each, where only one rerouting happens to the longest flow.

2. **Rank 4 (*full-rank*)**. The servers are not grouped in any way. All servers participate to the priority-driven load balance. Therefore, there are 3 thresholds and the longest flow it is rerouted three times.

In the next experiments (Fig.4.7) is evaluated the full-rank spatial diversity system for three topologies, corresponding to  $K = 3$ ,  $K = 5$ ,  $K = 9$  and for two serving policies: PS and FIFO. The number of priority queues  $N$  is still one per server. Remember that the numerical simulator is a job level simulator, thus flows are not fragmented in packets anyhow. This means that the FIFO policy is absolutely non-preemptive for jobs in the same priority queue. The service of a flow cannot be interrupted and flows arriving in the meanwhile are queued back in a first-in first-out order. Instead, the PS discipline subdivides the available capacity among all flows present in the system with a fluid approximation. Thus, a fresh flow share immediately the service rate granted to its priority queue with other flows in the same queue. These experiments were

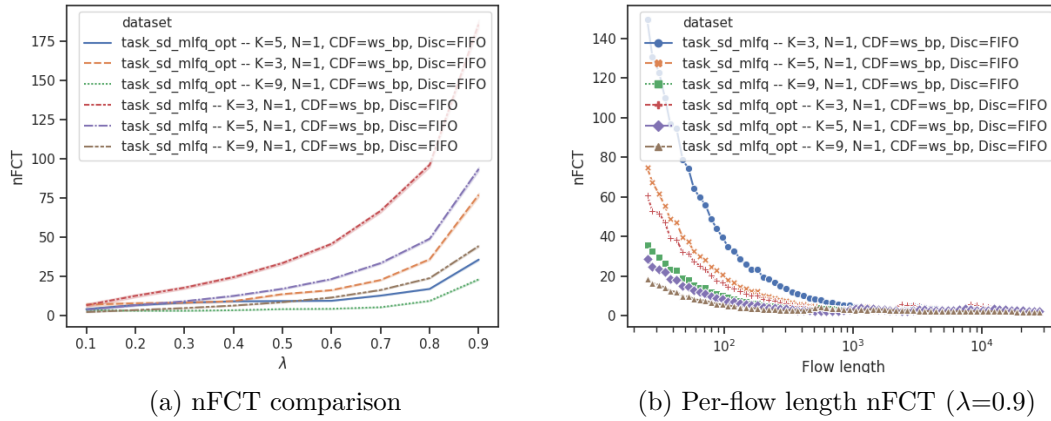


Figure 4.7: Web Search workload and FIFO discipline at 99% confidence interval.

carried out both with and without the optimized load balance thresholds. The first very positive observation is that, whatever the number of servers, the optimal load balance wins over the proportionate split. This again confirms the validity of the queuing model formulation for load balancing. The FIFO policy is clearly unfair with respect to short flows (Fig.4.7b). This is because without preemption a single elephant flow could starve a myriad of short flows. The unfairness is mitigated when increasing the spatial diversity rank, as long flows are demoted earlier and leave quickly high priority servers. Notably, for a fixed  $K$ , having  $N=1$  priority queue is the worst case from the point of view of mice flow starvation. With  $N>1$  PQs, the service received by a flow on each server

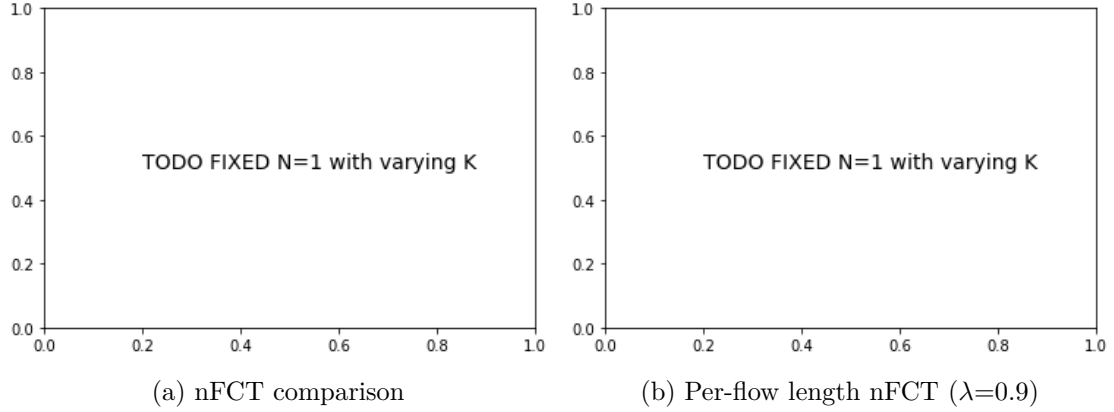


Figure 4.8: Web Search workload and PS discipline at 99% confidence interval. **GRAFICI MANCANTI MA SIMULAZIONI GERMAN CON QUELLO SCENARIO GIA FATTE MA USATE NEI GRAFICI FINALI. CHIEDERE GRAFICO A PARTE SE UTILE**

would be broken in multiple phases, each on a different priority queue. Since the PQs are scheduled in Strict Priority order and newly arrived flows enter in high priority, a short flow behind longer ones would have to wait only until the demotion at lower priority of the flows ahead. Somehow the demotion preempts the long flows, and weakens the impact of having a FIFO policy on the PQs themselves. The scenario changes with the PS service discipline. Short flows do not suffer the same issues yet mentioned, because they always share the processor with other flows. However, when increasing the rank long flows start to be excessively penalized. Remember that subsequent servers observe cut workloads, whose heavy-tailed property is gradually destroyed. In other words, the flow size distributions fed in low priority servers have reduced variability (Fig. 4.1a). This is a clear downside for PS policy. Indeed, it is well-known that under the assumption of low variability of the distribution, it is better to adopt a FCFS order. As a trivial example, suppose to schedule two flows  $f_1, f_2$  of size 10 bits/bytes/packets,... Denote as  $T(f_i)$  the average flow completion time of flow  $i$ , expressed in transmission slots. FIFO scheduling would give:

$$T_{FIFO}(f_1) = T_{FIFO}(f_2) = \frac{10 + 20}{2} = 15$$

Instead, PS scheduling would be worse:

$$T_{PS}(f_1) = T_{PS}(f_2) = \frac{19 + 20}{2} = 19.5$$

**TODO PRECEDENTE DISCORSO VA VERIFICATO CON IL GRAFICO SOPRA MANCANTE**

The study of the system with increased number of priority queues led to even more surprising results. This step really integrates the spatial diversity in the legacy MLFQ system with many priority levels. Here it is reported only the Processor Sharing discipline for the reasons explained few lines above. It is clear that for the FIFO case increasing  $N$  further minimize the average flow completion time. A simulative result confirming this

fact will be present in the last section §?? of the chapter during the final comparisons of SD-MLFQ (spatial diversity case) with ES-N (non spatial diversity case). Figure 4.9 is a key result that will uncover potential bottlenecks of SD-MLFQ. In this experiments are compared the average flow completion times of the system with fixed  $K=4$  and variable  $N$ . Load balance thresholds are not optimized in this case, but are set to grant the proportionate split of the traffic on the four parallel servers. The sub-thresholds have been assigned with the simple ES-N variation for spatial diversity (§3.2.1). Very interestingly, this numerical results provide a different perspective on the ideal granularity of priority queues per server with respect to previous works. While on the systems that approximate LAS or SRPT schedulers with PQs [3, 4], increasing the number of priority queues is always beneficial, it is not so when spatial diversity is introduced. At low load there is not a visible difference among the simulated scenarios. However, at high load the best average nFCT is attained with  $N=2$  priorities per server (dashed orange curve). The peculiar phenomenon is the behavior of the case  $N=8$  (red curve), whose corresponding nFCT starts to increase consistently above  $\lambda = 0.8$ . From the detailed breakdown of Fig.4.9b,

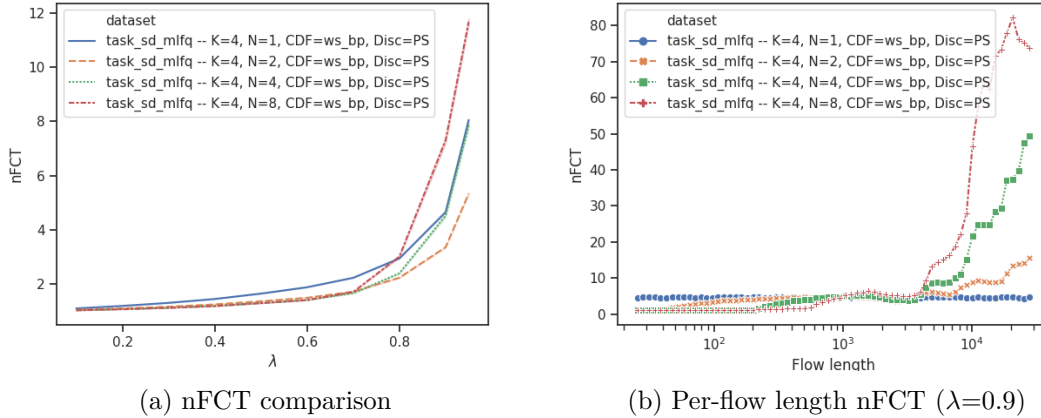


Figure 4.9: Web search workload and PS discipline at 99% confidence interval

it is evinced that the most severe degradation happens to long flows. They suffer an unacceptable penalty, which globally makes also the average worse. As a general trend, an higher number of priorities per server is reflected in a more pronounced raise of the job response time of long flows. At the same time, coherently with the scheduling theory of LAS (§2.1.2) and the results of PIAS, flows between 100KB and 1MB (medium size flows) are improved.

Such an extreme behavior on long flows was worth of extra-attention. In particular, it brings up an important question: is the best choice to configure low priority servers with a fine-grained prioritization?

### **Elephant demotion**

TODO Describe figures of hazard rate and optimal demotions for low priority spines. Link this observation with previous intuition that subsequent spines observe distributions with less and less variability (LAS and SRPT require high variability !!!).

### **Flow synchronization**

Problems of previous subsection are exacerbated by flow syn phenomenon (not observed in PIAS).

### **4.2.2 Final comparison**

Final comparison with ES-N (no spatial). Show plots and resume table.





## Chapter 5

# Evaluation in a Datacenter Network

### 5.1 title

#### 5.1.1 title

#### 5.1.2 title

### 5.2 title

#### 5.2.1 title

#### 5.2.2 title



# Bibliography

- [1] Samuli Aalto, Urtzi Ayesta, and Rhonda Righter. On the gittins index in the m/g/1 queue. *Queueing Systems*, 63(1-4):437–458, 2009.
- [2] Mohammad Alizadeh, Albert Greenberg, David A. Maltz, Jitendra Padhye, Parveen Patel, Balaji Prabhakar, Sudipta Sengupta, and Murari Sridharan. Data center tcp (dctcp). In *Proceedings of the ACM SIGCOMM 2010 Conference*, SIGCOMM '10, page 63–74, New York, NY, USA, 2010. Association for Computing Machinery.
- [3] Mohammad Alizadeh, Shuang Yang, Milad Sharif, Sachin Katti, Nick Mckeown, Balaji Prabhakar, and Scott Shenker. pfabric. *Proceedings of the ACM SIGCOMM 2013 conference on SIGCOMM - SIGCOMM 13*, 2013.
- [4] Wei Bai, Li Chen, Kai Chen, Dongsu Han, Chen Tian, and Hao Wang. Pias: Practical information-agnostic flow scheduling for commodity data centers. *IEEE/ACM Transactions on Networking*, 25(4):1954–1967, 2017.
- [5] Wei Bai, Li Chen, Kai Chen, and Haitao Wu. Enabling ecn in multi-service multi-queue data centers. In *Proceedings of the 13th Usenix Conference on Networked Systems Design and Implementation*, NSDI'16, page 537–549, USA, 2016. USENIX Association.
- [6] Paul J. Burke. The output of a queuing system. *Operations Research*, 4(6):699–704, 1956.
- [7] Mark E. Crovella. Performance evaluation with heavy tailed distributions. In Dror G. Feitelson and Larry Rudolph, editors, *Job Scheduling Strategies for Parallel Processing*, pages 1–10, Berlin, Heidelberg, 2001. Springer Berlin Heidelberg.
- [8] Albert Greenberg, James R. Hamilton, Navendu Jain, Srikanth Kandula, Changhoon Kim, Parantap Lahiri, David A. Maltz, Parveen Patel, and Sudipta Sengupta. Vl2: A scalable and flexible data center network. In *Proceedings of the ACM SIGCOMM 2009 Conference on Data Communication*, SIGCOMM '09, page 51–62, New York, NY, USA, 2009. Association for Computing Machinery.
- [9] Chuanxiong Guo, Haitao Wu, Zhong Deng, Gaurav Soni, Jianxi Ye, Jitu Padhye, and Marina Lipshteyn. Rdma over commodity ethernet at scale. In *Proceedings of the 2016 ACM SIGCOMM Conference*, SIGCOMM '16, page 202–215, New York, NY, USA, 2016. Association for Computing Machinery.
- [10] J. Kennedy and R. Eberhart. Particle swarm optimization. In *Proceedings of ICNN'95 - International Conference on Neural Networks*, volume 4, pages 1942–1948 vol.4, 1995.
- [11] Idris A. Rai, Guillaume Urvoy-Keller, and Ernst W. Biersack. Analysis of las scheduling for job size distributions with high variance. *Proceedings of the 2003 ACM*

- SIGMETRICS international conference on Measurement and modeling of computer systems - SIGMETRICS 03*, 2003.
- [12] Arjun Roy, Hongyi Zeng, Jasmeet Bagga, George Porter, and Alex C. Snoeren. Inside the social networks (datacenter) network. *Proceedings of the 2015 ACM Conference on Special Interest Group on Data Communication - SIGCOMM 15*, 2015.
  - [13] Linus Schrage. A proof of the optimality of the shortest remaining processing time discipline. *Operations Research*, 16(3):687–690, 1968.
  - [14] Y. Shi and R. Eberhart. A modified particle swarm optimizer. In *1998 IEEE International Conference on Evolutionary Computation Proceedings. IEEE World Congress on Computational Intelligence (Cat. No.98TH8360)*, pages 69–73, 1998.
  - [15] David J Wales and Jonathan PK Doye. Global optimization by basin-hopping and the lowest energy structures of lennard-jones clusters containing up to 110 atoms. *The Journal of Physical Chemistry A*, 101(28):5111–5116, 1997.









