



# Designing Data Center Networks Using Bottleneck Structures

Jordi Ros-Giralt<sup>1</sup>, Noah Amsel<sup>1</sup>, Sruthi Yellamraju<sup>1</sup>, James Ezick<sup>1</sup>, Richard Lethin<sup>1</sup>, Yuang Jiang<sup>2</sup>, Aosong Feng<sup>2</sup>,  
Leandros Tassioulas<sup>2</sup>, Zhenguo Wu<sup>3</sup>, Min Yee Teh<sup>3</sup>, Keren Bergman<sup>3</sup>  
Reservoir Labs<sup>1</sup>, Yale Institute for Network Science<sup>2</sup>, Columbia University<sup>3</sup>  
New York City, NY, USA<sup>1</sup> New Haven, CT, USA<sup>2</sup> New York City, NY, USA<sup>3</sup>  
giralt@reservoir.com

## ABSTRACT

This paper provides a mathematical model of data center performance based on the recently introduced Quantitative Theory of Bottleneck Structures (QTBS). Using the model, we prove that if the traffic pattern is *interference-free*, there exists a unique optimal design that both minimizes maximum flow completion time and yields maximal system-wide throughput. We show that interference-free patterns correspond to the important set of patterns that display data locality properties and use these theoretical insights to study three widely used interconnects—fat-trees, folded-Clos and dragonfly topologies. We derive equations that describe the optimal design for each interconnect as a function of the traffic pattern. Our model predicts, for example, that a 3-level folded-Clos interconnect with radix 24 that routes 10% of the traffic through the spine links can reduce the number of switches and cabling at the core layer by 25% without any performance penalty. We present experiments using production TCP/IP code to empirically validate the results and provide tables for network designers to identify optimal designs as a function of the size of the interconnect and traffic pattern.

## CCS CONCEPTS

• **Networks** → **Network performance modeling; Data center networks; Network design principles;**

## KEYWORDS

Data center, design, model, bottleneck structure, max-min

---

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

SIGCOMM '21, August 23–27, 2021, Virtual Event, USA

© 2021 Copyright held by the owner/author(s).

ACM ISBN 978-1-4503-8383-7/21/08.

<https://doi.org/10.1145/3452296.3472898>

## 1 INTRODUCTION

Data centers are some of the largest centrally-managed networks in the world, responsible for storing, computing and distributing large amounts of data. Driven by ever increasing machine-to-machine workloads [28], network researchers and architects have focused on identifying data center topologies and designs that are able to scale and deliver high-performance at low cost [1, 18, 28, 30]. The steep cost of building and operating these interconnects motivates data centers to employ techniques such as oversubscription and bandwidth tapering [1, 21], which allow them to optimize the performance-cost trade-off. However, because little is known about the mathematical principles that drive the performance of data center interconnects, network architects are forced to either take conservative approaches that waste bandwidth and unnecessarily increase costs or to use trial-and-error methodologies that are operationally very costly.

In this paper, we position a formal mathematical model of data center networks. Our approach is analogous to the way physicists study natural systems—by formulating a set of fundamental laws and using them to derive specific predictions about the system’s behavior. We begin from the assumption that the network is congestion-controlled.<sup>1</sup> Congestion control algorithms are designed to regulate traffic according to the following principle: *maximize the throughput of each flow while ensuring fairness among flows*. This principle ensures that the resources of the network are not wasted (which would not maximize throughput) and that no flow is entirely starved of the bandwidth it needs (which would be unfair).

The recently discovered Quantitative Theory of Bottleneck Structures (QTBS) [26, 27] builds on this assumption to create a mathematical model of congestion-controlled networks. In this paper, we use it to study the problem of designing data center networks. Our model predicts how a data center network will perform for a given traffic pattern. This allows us to derive optimal network designs and identify wasteful ones using our model alone, without resorting to costly trial-and-error deployments. For instance, the model reveals that a 3-level folded-Clos interconnect [1] with radix 24 that routes

---

<sup>1</sup>Congestion control is implemented as part of the TCP and Infiniband protocols used in many data center networks.

10% of the traffic through the spine links can reduce the number of switches and cabling at the core layer by 25% without any performance penalty.

The contributions of this paper are as follows:

- We introduce a class of network designs called *proportional designs*. We show that for traffic patterns that we call *interference-free*, proportional designs achieve maximal throughput and minimal flow completion time for any fixed cost (Section 2).
- We show that three widely-used types of data center topologies—fat-tree (Section 3), folded-Clos (Section 4), and dragonfly (Appendix G)—are interference free for common traffic patterns that exhibit data locality properties.
- We derive optimal designs for these topologies and identify designs that are wasteful and ought to be avoided (Sections 3 and 4).
- Through extensive simulations using production-grade TCP/IP code, we demonstrate that the predictions of our theory hold for the three studied topologies (Section 5).

## 2 THEORETICAL BACKGROUND

### 2.1 Network Model

In our network model, a data center consists of a set of hosts that are interconnected using a set of links. Each link has a finite capacity and hosts communicate with each other using flows. A flow between a pair of hosts traverses a subset of the links, which we call its path. We assume the network is regulated by a congestion control algorithm that determines the transmission rate of each flow according to the following principle: maximize the throughput of each flow while ensuring fairness among flows. More specifically, our mathematical model assumes the congestion control algorithm determines the rate of each flow according to the well-known max-min fairness criterion [6, 8]. While we acknowledge that production-grade congestion control algorithms do not precisely act according to max-min fairness, recent work [26, 27] demonstrates that the max-min assumption enables a powerful mathematical model of the network that can be used to analyze and predict network performance.

Our work builds on the quantitative theory of bottleneck structures (QTBS) introduced in [26, 27]. QTBS provides a general-purpose mathematical framework to model communication networks. Because of its predictive power, QTBS can be used to address a wide variety of communication problems such as traffic engineering, routing, flow scheduling, network design, capacity planning, resiliency analysis, network slicing, or service level agreement (SLA) management, among others, by incorporating an understanding of congestion control into its solution. In this paper, we use the assumptions of the QTBS model to develop a formal framework to design data center interconnects. Note that, while the reader is encouraged to review the existing literature on

QTBS [26, 27], the present paper is self-contained and does not require prior knowledge of it.

### 2.2 Designing Data Center Networks

In our framework, network architects are given a class of topology (e.g., folded-Clos) and the required size of the interconnect (i.e., the number of hosts to be connected). Their objective is to set the capacity of each link so that the interconnect achieves maximal performance at the lowest possible cost. We refer to an assignment of capacity values to the links of an interconnect as a *design*:

*Definition 2.1. Network design.* Let  $\mathcal{L}$  be the set of links in a given data center network. A *network design* (or simply a *design*) is a function  $c : \mathcal{L} \rightarrow \mathbb{R}_{>0}$  mapping each link to its capacity. We use the notation  $c_l$  interchangeably with  $c(l)$ .

Careful network design is increasingly important, as improper allocation of bandwidth in an interconnect can lead to wasteful capital and operational expenditures at the scale of modern data centers [28, 30]. In performing our analysis, we need two kinds of inputs: information about the topological structure of the network and assumptions about the traffic pattern that it needs to support. We describe the topologies we study in detail in Sections 3 (fat-trees), 4 (folded-Clos) and Appendix G (dragonflies), where the mathematical framework for each specific interconnect is presented. As for the traffic pattern, we characterize it as follows:

*Definition 2.2. Traffic Pattern.* Let  $\mathcal{H}$  be the set of hosts in a given interconnect. Then a *traffic pattern* is a function  $b : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}_{\geq 0}$  mapping each ordered pair of hosts to the amount of data (e.g., in bits) to be transmitted from the first to the second.

We say that the traffic pattern is *uniform* if all pairs of hosts transmit the same number of bits, and *skewed* otherwise. We use the shorthand  $\mathcal{F} = \mathcal{H} \times \mathcal{H}$  to denote the set of data flows connecting every pair of hosts in the network. We assume that for each ordered pair of hosts  $h, h' \in \mathcal{H}$ , there exists a single data flow  $f$  transmitting  $b(f)$  data from  $h$  to  $h'$ . We can now introduce our criteria for evaluating a network's performance:

*Definition 2.3. Network completion time.* Let  $\mathcal{L}$  and  $\mathcal{F}$  be the sets of links and flows, and let  $b$  be the traffic pattern. Assume all flows start transmitting data at the same time. For a given design  $c$ , let  $\text{fct}(b, c, f)$  be the time it takes for flow  $f$  to complete transmitting its  $b(f)$  units of data. Then,  $\mu(b, c) = \max_{f \in \mathcal{F}} \text{fct}(b, c, f)$  is the completion time of the network for the given traffic pattern and design.

Since network completion time refers to the time it takes to complete the longest duration flow, minimizing network completion time provides a mechanism for reducing flow *completion time tail*. The latter is widely regarded as one of the key performance metrics in data centers [2, 11, 23, 32],

because a large flow completion time tail causes latency and jitter, harming application performance. In this paper, we show that the optimal designs derived from our framework equalize the completion time of all flows (see Theorem 2.10), thus reducing to zero the tail of the flow completion time distribution. Moreover, network completion time is closely connected to network throughput, another key metric of system-wide performance, as follows. Assume that instead of just one flow  $f$  between each pair of hosts, the traffic pattern now contains  $n$  identical copies of  $f$  (or  $n$  copies of the original batch), each of which transmits the same amount of data  $b(f)$ . Unlike before, flows can be scheduled arbitrarily, and flows from different batches may be transmitting simultaneously. The throughput of a network measures the average rate at which it can transmit data as the number of batches grows:

*Definition 2.4. Network throughput.* For a given traffic pattern  $b$  and design  $c$ , let  $\text{bct}(b, c, n)$  be the smallest possible time to complete transmitting  $n$  batches when using the best scheduling. Then, the *network throughput* of the design is  $T(b, c) = \lim_{n \rightarrow \infty} n/\text{bct}(b, c, n)$ .

Note that if the network completion time is  $\mu(b, c)$ , the throughput is at least  $1/\mu(b, c)$ , since one can simply schedule the batches sequentially. In some cases, it is possible to achieve higher throughput by scheduling batches to partially overlap. However, this scheduling harms latency; the first batch takes longer to complete, since it must now share resources with the second batch. Designers must balance the goals of achieving high performance and low cost. Requirements differ from case to case, and so the best design depends on context. However, some designs are simply wasteful in that their cost is unnecessarily high for the level of performance they achieve. Ideally, a network is designed so that all of its capacity is utilized throughout the transmission of the traffic pattern. Intuitively, this means that the resources of the network are being used efficiently, as the operator is not paying for bandwidth that is idle. In this case, scheduling successive batches sequentially achieves optimal throughput without sacrificing latency. Since no bandwidth is wasted during transmission of the first batch, it is impossible to increase throughput by packing in flows from the second batch. The following definition and theorem capture this intuition and relate it to the criteria discussed above:

*Definition 2.5. Wasteless designs.* For a given topology and traffic pattern  $b$ , a design  $c$  is wasteless if all the bandwidth of each link is used throughout the transmission. That is,  $\forall t \in [0, \mu(b, c)]$ ,  $\sum_{f \in \mathcal{F}_l} r_{c,b}(f, t) = c(l)$ , where  $\mathcal{F}_l$  is the set of flows that traverse  $l$  and  $r_{c,b}(f, t)$  is the rate of flow  $f$  at time  $t$  for design  $c$  and traffic pattern  $b$  according to the congestion control algorithm.

**THEOREM 2.6. Optimality of wasteless designs.** For a given topology and traffic pattern, if a design  $c$  is wasteless, then it is impossible to improve on the network completion time or network throughput of the design without adding more capacity. That is, if  $c'$  is an alternative design for which  $\mu(b, c') < \mu(b, c)$  or  $T(b, c') > T(b, c)$ , then  $\sum_{l \in \mathcal{L}} c'(l) > \sum_{l \in \mathcal{L}} c(l)$ .

PROOF. See Appendix A.1 □

Intuitively, for a given topology, a wasteless design achieves the best possible performance for a fixed cost, where cost is a function of the total capacity of the network. Note that in production networks, operators typically do not design interconnects to be wasteless. Rather, networks are provisioned with extra capacity to accommodate for potential link failures and to protect latency-sensitive flows that can arrive at arbitrary times. However, network designers must know how to provision their networks before they can know how to overprovision them. That is, they need to first identify the optimal design they would use if link failures, latency, and traffic variability did not exist, and then provision excess capacity to accommodate for these factors on top of the base design. Identifying the wasteless design first gives network designers a principled, quantitative way to tell how much they are overprovisioning each link. Our framework aims to help operators identify this base design.

### 2.3 Proportional Designs

To address the data center design problem described above, we introduce the following class of designs:

*Definition 2.7. Proportional Designs.* For a given traffic pattern  $b$ , a design  $c$  is *proportional* if each link's capacity is proportional to the sum of the sizes of the flows that traverse it. That is,

$$c(l) = \alpha \sum_{f \in \mathcal{F}_l} b(f)$$

for some  $\alpha > 0$ , where  $\mathcal{F}_l$  is the set of flows that traverse  $l$ .

A proportional design exists for any topology and traffic pattern and is unique up to the scaling factor  $\alpha$ . By choosing  $\alpha$ , one can create a proportional design that achieves any desired level of performance. Our main theoretical result motivates the use of proportional designs: proportional designs are optimal in the sense that no other design can achieve zero waste, as shown next.

**THEOREM 2.8. Non-proportional designs waste bandwidth.** If a design wastes no bandwidth on a given traffic pattern, then it is the proportional design for that traffic pattern.

PROOF. See Appendix A.2. □

The converse of this theorem is not true in general, because for some traffic patterns, it is impossible not to waste bandwidth even if a proportional design is used. (For an example, see Lemma 3.3.) The following definition and theorem

give a sufficient condition on the traffic pattern to ensure that a wasteless design exists (and, by extension, that the proportional design is wasteless):

**Definition 2.9.** *Interference-free.* For a given topology, a traffic pattern  $b$  is said to be *interference-free* if each flow  $f$  traverses some link  $l$  that is traversed by no flow that transmits more bits than  $f$ . That is,  $\forall f \in \mathcal{F}, \exists l \in \mathcal{L}$  s.t.  $f \in \arg \max_{f' \in \mathcal{F}_l} b(f')$ .

**THEOREM 2.10.** *Interference-free proportional designs are wasteless. If the traffic pattern is interference free, then the proportional design wastes no bandwidth and all flows finish transmitting at the same time.*

**PROOF.** See Appendix A.3.  $\square$

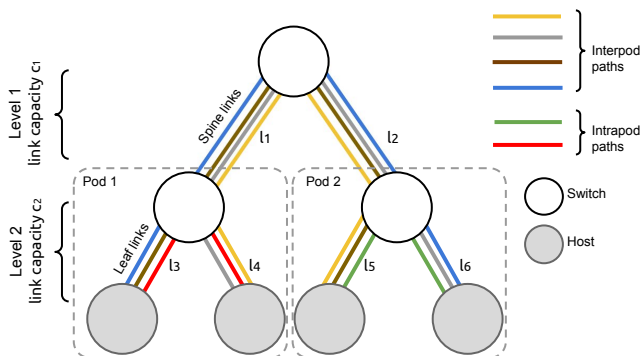
To sum up, the above theorem shows that if a traffic pattern is interference free, then proportional designs (and only proportional designs) have all of the following optimality properties: (1) no bandwidth is wasted; (2) all flows terminate at the same time; (3) the network completion time is the smallest possible without adding more capacity to the network; (4) the network throughput is as large as possible without adding more capacity to the network. Further, in the following sections, we show that typical data center traffic patterns that display locality properties are interference-free, ensuring that proportional designs are optimal in these ways.

### 3 DESIGNING FAT-TREE NETWORKS

Fat-trees are a popular class of topologies first introduced by Leiserson. As he demonstrated, fat-trees are universally efficient networks in the following sense: for a given network size  $s$ , a fat-tree can emulate any other network of that size  $s$  with a performance slowdown at most logarithmic in  $s$  [20]. This property makes fat-tree topologies highly competitive and is one of the reasons they are widely used in large-scale data centers [1] and high-performance computing (HPC) networks [21]. A fat-tree is a complete  $n$ -ary tree [9] that satisfies the condition  $c_l > c_{l'}$  for pairs of links such that  $l$  is nearer to the root than  $l'$  [20]. Fat-trees are said to be *full* if every level has the same total capacity, that is:

$$\sum_{l \in \mathcal{L}_i} c_l = \sum_{l \in \mathcal{L}_j} c_l, \text{ for } 1 \leq i, j \leq L, \quad (1)$$

where  $\mathcal{L}_i$  is the set of links at level  $i$  of the tree and  $L$  is the tree's total number of levels [10]. We will use the notation  $FT(n, L)$  to denote a fat-tree in which each switch has  $n$  children and the total number of levels is  $L$ . (Note that  $L$  is half the diameter of the fat-tree.) For example, Fig. 1 shows the topology of a  $FT(2, 2)$ . We will also consider a generalization of fat-trees in which nodes at different levels of the tree may have different numbers of children, so long as all nodes at a single given level have the same number of children. We use the notation  $FT([n_1, n_2, \dots, n_L])$  to denote such a tree, where each switch at level  $i$  of the tree has  $n_i$  children—thus,



**Figure 1: A fat-tree with two levels and two children per switch, denoted as  $FT(2, 2)$ .**

$FT(n, L) = FT([n, n, \dots, n])$ . Note that because there is only a single path between every pair of hosts, routing is trivial.

For the specific case of fat-trees with 2-levels ( $L = 2$ ), we adopt the common terminology of *spine* and *leaf* links to refer to the upper-level links (connected to the root) and the lower-level links (connected to the hosts), respectively. We also use the terms *interpod* and *intrapod* paths to refer to paths that traverse the spine links and paths that traverse only the leaf links, respectively. Because every host communicates with every other host, each path accommodates two flows, one in each direction. For instance, the topology  $FT(2, 2)$ , shown in Fig. 1, has four hosts, two interpod links, four intrapod links, four interpod paths, two intrapod paths, eight interpod flows and four intrapod flows.

In our analysis, we will assume all links at level  $i$  of a fat-tree have the same capacity value  $c_i$ , as is the case in most production deployments [1, 30]. For the widely used case of 2-level  $n$ -ary fat-trees ( $L = 2$ ), we refer to the design parameter  $\tau = c_1 / (n \cdot c_2)$  as the *tapering parameter* of the network. It is easy to see from Equation (1) that  $\tau = 1$  is the case of a full fat-tree [10], providing the full bisectional bandwidth in all fat-tree levels. In general, this parameter is the ratio of the actual spine link capacity divided by the spine link capacity of a full fat-tree network with the same leaf link capacity as the given network. The tapering parameter  $\tau$  determines the performance/cost trade-off of the interconnect: increasing  $\tau$  improves performance but it also increases the cost of the interconnect, and vice versa. This parameter also characterizes the degree to which the spine links are oversubscribed, usually expressed with the notation  $1/\tau:1$  [1]. For instance, an  $FT(n, 2)$  design with a tapering parameter of 0.5 is oversubscribed by a factor of 2:1.

#### 3.1 Design Equations for Uniform Traffic

We start by solving the fat-tree design problem under the assumption of uniform traffic:

LEMMA 3.1. *Optimal fat-tree with uniform traffic. Consider a generalized fat-tree  $FT([n_1, n_2, \dots, n_L])$  and let  $c_i$  be the capacity of its links at level  $i$ , for  $1 \leq i \leq L$ . Then, if the traffic pattern is uniform, it is interference-free and the following design is optimal:*

$$c_i = \frac{\rho_i}{\rho_L} c_L \quad (2)$$

where

$$\rho_i = 2 \left( \prod_{j=1}^i n_j - 1 \right) \prod_{j=i+1}^L n_j^2 \quad (3)$$

PROOF. See Appendix A.5.  $\square$

This equation comes from applying the definition of a proportional design to a fat-tree with uniform traffic. Furthermore, since the uniform traffic pattern is interference-free, Theorem 2.10 states that this design equalizes the completion time of all the flows. Under the assumption of uniform traffic where all hosts send the same amount of data to each other, this is equivalent to equalizing the throughput of all the flows.

For example, the set of optimal two-level fat-trees  $FT(n, 2)$  can be derived by using Equations (2) and (3) while setting  $L = 2$  and  $n_1 = n_2 = n$ , where  $n$  is the number of children of each switch in the tree. Doing some simple algebraic manipulations, this leads to the following optimal design:

$$c_1 = \frac{n^2}{n+1} c_2 \quad (4)$$

Fig. 2 represents the above designs using a red line. For a fixed leaf link capacity, the capacity of spine links increases quasi-linearly. Of special interest is also the set of designs corresponding to the full fat-tree solution, represented by the equation  $c_1 = n \cdot c_2$  and shown in Fig. 2 as a black line. We can see that a full fat-tree is not an optimal design when traffic is uniform, because some of the bandwidth of the spine links goes to waste. In other words, the network operator is paying for bandwidth in the spine links that is never used. In the next section, we generalize this analysis of fat trees to skewed (non-uniform) traffic patterns.

### 3.2 Design Equations for Skewed Traffic

We now introduce the general equations of an optimal fat-tree design for skewed traffic:

LEMMA 3.2. *Optimal fat-tree with skewed traffic. Consider a generalized fat-tree  $FT([n_1, n_2, \dots, n_L])$  and let  $c_i$  be the capacity of its links at level  $i$ , for  $1 \leq i \leq L$ . Assume a traffic pattern  $b(f) = \sigma_i$ , where  $f$  is a flow that traverses a link in level  $i$ , but no link in level  $i-1$ .<sup>2</sup> If  $i \leq j \implies \sigma_i \leq \sigma_j$  for all  $1 \leq i, j \leq L$ , then the traffic pattern is interference free and*

<sup>2</sup>So a flow that goes through all levels of the tree transmits  $\sigma_1$ , and a flow between a pair of hosts that are nearest neighbors in the tree transmits  $\sigma_n$ .

the following design is optimal:

$$c_1 = \pi_1 \sigma_1 \quad (5)$$

$$c_i = \pi_i \sigma_i + \frac{c_{i-1}}{n_i} \quad (6)$$

where

$$\pi_i = (n_i - 1) \prod_{j=i+1}^L n_j^2 \quad (7)$$

PROOF. See appendix A.6.  $\square$

In addition to providing the equations of the optimal design, the above lemma characterizes the set of interference-free traffic patterns in a fat-tree topology. In particular, a traffic pattern is interference free if  $i \leq j \implies \sigma_i \leq \sigma_j$  for all  $1 \leq i, j \leq L$ . This corresponds precisely to the set of traffic patterns that display locality properties—patterns where hosts transmit more data to other hosts in their same pod than to hosts in other pods. Because good design principles dictate that operators should map applications onto a data center by exploiting data locality, from a design standpoint, we reason that interference-freeness is the property that characterizes the set of interesting traffic patterns to achieve best performance in data center networks.

Applying some algebraic manipulations on Equations (5), (6) and (7), we can derive the set of optimal designs for a 2-level fat-tree  $FT(n, 2)$  by setting  $L = 2$ ,  $n_1 = n_2 = n$ , and letting  $\sigma = \sigma_2/\sigma_1$ , which leads us to:

$$c_1 = \frac{n^2}{n+\sigma} c_2 \quad (8)$$

(Note that the above equation corresponds to Equation (4) when traffic is uniform.) This leads to the following optimal tapering parameter:

$$\tau = \frac{c_1}{n \cdot c_2} = \frac{n}{n+\sigma} \quad (9)$$

Or, equivalently, an oversubscription of  $(n+\sigma)/n:1$ . For instance, a  $FT(16, 2)$  in which intrapod flows carry ten times more traffic than interpod flows ( $\sigma = 10$ ) has an optimal tapering parameter  $\tau = 0.61538$  (an oversubscription of 1.625:1), yielding the design  $c_1/c_2 = 9.8462$ . This design is shown in Fig. 2 as a blue dot. Similarly, the same interconnect with uniform traffic ( $\sigma = 1$ ) has an optimal tapering parameter  $\tau = 0.94118$  (an oversubscription of 1.0625:1), yielding the design  $c_1/c_2 = 15.0588$ . This design is shown in Fig. 2 as a green dot.

The following lemma characterizes the design space of fat-trees  $FT(n, 2)$ :<sup>3</sup>

LEMMA 3.3. *Design space of fat-trees  $FT(n, 2)$ . Consider a fat-tree  $FT(n, 2)$  and let  $c_1$  and  $c_2$  be the capacity of its spine and leaf links, respectively. Without loss of generality, let*

<sup>3</sup>This lemma can be generalized to support  $FT([n_1, n_2, \dots, n_L])$ , we leave this as an exercise to the reader.



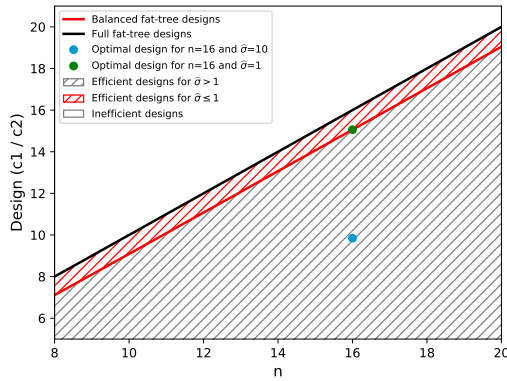


Figure 2: Design space for  $FT(n, 2)$ .

$b(f) = 1$  for interpod flows, and  $b(f) = \sigma$  for intrapod flows. Then,

- (1) If  $\sigma \geq 1$ , the optimal design satisfies  $c_1 \leq n^2 \cdot c_2 / (n + 1)$ .
- (2) If  $0 < \sigma < 1$ , the proportional design satisfies  $n^2 \cdot c_2 / (n + 1) < c_1 < n \cdot c_2$ , but every design wastes bandwidth.
- (3) If  $\sigma = 0$ , the optimal design corresponds to  $c_1 = n \cdot c_2$ .

PROOF. See appendix A.7. □

The above lemma is pictured in Fig. 2 as follows. The region of optimal designs for traffic patterns such that  $\sigma \geq 1$  (traffic patterns that are interference-free) is marked with gray hash lines. It is delimited by the optimal design line for uniform ( $\sigma = 1$ ) traffic  $c_1 = n^2 \cdot c_2 / (n + 1)$ , shown as a red line. Note that the optimal designs for  $\sigma > 1$  correspond to oversubscribed interconnects. The region of optimal designs for  $\sigma < 1$  (traffic patterns that are not interference-free) is marked with red hash lines, bordered by the full fat-tree design line  $c_1 = n \cdot c_2$  and the optimal design line for uniform traffic. These designs correspond to undersubscribed interconnects. The design space lemma also shows that there exists no traffic pattern for which a design in the region above full fat-trees ( $c_1 > n \cdot c_2$ ) is optimal. That's because for any traffic pattern, such a design would require more bandwidth in the spine links than a full fat tree with the same leaf links would, but delivers exactly the same performance.

In practice, data centers tend to experience traffic skewness ( $\sigma > 1$ ), as applications exploit data locality by sharing more data with other hosts in their own pod than to hosts in other pods [4, 5, 17]. Thus, for typical traffic patterns, full fat-trees are inefficient interconnects too. Furthermore, a full fat-tree's inefficiency increases as  $\sigma$  increases. This result may appear surprising since Leiserson demonstrated that full fat-trees are universally efficient [20]. However, Leiserson's

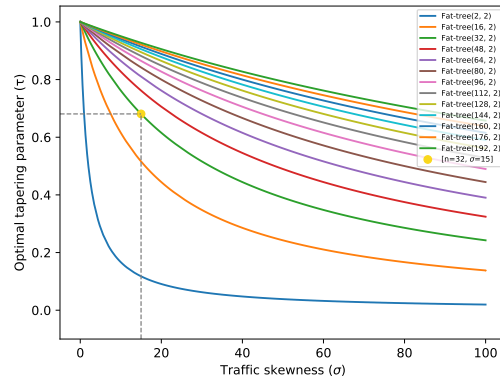


Figure 3: Optimal tapering parameter as a function of traffic skewness for various  $FT(n, 2)$  designs.

interconnect assumes a best-effort message-based communication system without congestion control regulation<sup>4</sup>. The presence of congestion control in modern interconnects has the effect of shifting the optimal design depending on the network's traffic pattern. While full fat-trees are rarely used in production networks due to their prohibitive costs [1], this result reaffirms that it is wasteful to do so, since the presence of traffic skewness implies that the optimal design is oversubscribed, with an oversubscription factor that depends on the degree of traffic skewness  $\sigma$  according to Lemma 3.2.

Fig. 3 shows the optimal tapering parameter  $\tau$  as a function of traffic skewness  $\sigma$  for a variety of 2-level fat-trees. Once again, we see that all optimal designs correspond to oversubscribed interconnects (i.e.,  $\tau < 1$ ) departing away from the full fat-tree design ( $\tau = 1$ ) as traffic skewness increases. We also note that as the size of the interconnect  $n$  increases, the optimal tapering parameter increases. For example, suppose that our goal is to design a fat-tree  $FT(32, 2)$  to transport a traffic pattern with skewness  $\sigma = 15$ . By using the chart in Fig. 3, we can identify the needed design (represented with a yellow dot) at the intersection of the green curve ( $FT(32, 2)$ ) with the line  $\sigma = 15$ , resulting in a tapering parameter of  $\tau = 0.68085$ .

Appendix D presents a chart that plots the traffic skewness value needed to optimally operate a  $FT(n, 2)$  as a function of the number of hosts  $n^2$  supported by the interconnect and for various tapering parameters  $\tau$ . Appendix J includes additional design tables for the more general class of fat-trees  $FT([n_1, n_2])$  and various  $\sigma$  values derived from the general equations in Lemma 3.2.

<sup>4</sup>This is similar to the way the Internet operated prior to the invention of the first congestion control algorithm. Jacobson published the first congestion control algorithm for TCP/IP networks [14] three years after Leiserson's work [20].

## 4 DESIGNING FOLDED-CLOS NETWORKS

Our next goal is to apply QTBS to derive the equations that drive the optimal design of folded-Clos networks. While often times folded-Clos are also referred as fat-trees [1], mathematically we treat them separately because they have topological differences that lead to different equations. Throughout this section, we will use the notation introduced by Al-Fares et al. in their influential paper [1] to describe and analyze the performance of this important class of interconnects.

We will assume the folded-Clos network is implemented using switches of radix  $k$  and identical link capacity that are interconnected using a tree structure with  $L$  levels as described in [1]. These two parameters,  $k$  and  $L$ , characterize a class of folded-Clos topologies denoted as  $Clos(k, L)$ . Consider as an example the interconnect  $Clos(4, 3)$  illustrated in Fig. 4. It consists of 20 switches with radix 4, 48 links, 4 pods and 16 hosts interconnected via 120 possible paths. Links are organized in three levels ( $L = 3$ ), which are commonly referred from top to bottom as the *spine*, the *aggregation* and the *leaf* (or *edge*) levels. Folded-Clos interconnects with three levels have three possible types of flows, as shown in Fig. 4: *interpod flows* (red color) traversing the spine links, *long intrapod flows* (blue color) traversing the aggregation links but not the spine links, and *short intrapod flows* (yellow color) only traversing leaf links. As shown in [1], a general  $Clos(k, 3)$  interconnect has  $5k^2/4$  switches of radix  $k$ ,  $3k^3/4$  links,  $k^3/4$  hosts grouped in  $k$  pods, and  $k^6/32 - k^3/8$  possible paths. Finally, we will assume flows are routed according to the algorithm presented in [1]. Note that this is done without loss of generality, as we could apply our methodology to other types of routing (see also Section 7).

Most production data centers introduce oversubscription in the spine links as a mechanism to reduce the cost of the interconnect. In this case, the oversubscription of a folded-Clos topology corresponds to the ratio between the total capacity of the leaf links across all pods and the total capacity of the spine links, denoted as  $\omega:1$ . Consider for instance the  $Clos(4, 3)$  topology in Fig. 4. Because the total capacity of the leaf and spine links is the same (each group has 16 links of equal capacity), the oversubscription parameter is 1:1 ( $\omega = 1$ ). Suppose instead that we remove the two most-left spine switches from the network (and their corresponding eight links). Such a configuration yields an oversubscription of 2:1 ( $\omega = 2$ ).<sup>5</sup>

Because folded-Clos are discrete topologies, not all possible oversubscription values are available to a network designer. In particular, for the class of  $Clos(k, 3)$  interconnects,

<sup>5</sup>Note that the oversubscription parameter  $\omega$  is equal to the inverse of the bandwidth tapering parameter  $\tau$  introduced in Section 3. The main difference is that  $\omega \in \mathbb{Z}_{>0}$  while  $\tau \in \mathbb{R}_{>0}$ , reflecting the discrete nature of the folded-Clos interconnect.

oversubscription is implemented by grouping spine switches into blocks—called *spine blocks*—consisting of  $k/2$  switches each. Since there are  $k^2/4$  spine switches, this leads to a total of  $k/2$  possible oversubscription configurations of the form  $\omega:1$ , with  $\omega = \frac{k}{2\beta}$  and  $1 \leq \beta \leq k/2$ . Note that here  $\beta$  corresponds to the number of spine blocks deployed. With this approach, network architects can decide to increase the number of spine blocks in order to increase performance of the interconnect or to decrease the number of spine blocks to reduce its cost, providing an efficient and elegant mechanism to control its performance-cost trade-off. For instance, for a folded-Clos with radix  $k = 4$ , we have that  $\omega \in \{1, 2\}$  (since  $\beta \in \{1, 2\}$ ), which leads to two possible configurations, as shown in Fig. 4: 1:1 deploying both spine blocks (for a total of 4 spine switches) and 2:1 deploying only one of the spine blocks (for a total of 2 spine switches).

### 4.1 Design with Optimal Oversubscription

We now introduce the general equations that determine an optimal 3-level folded-Clos design for skewed traffic:

LEMMA 4.1. *Optimal 3-level folded-Clos with radix  $k$  and skewed traffic. Assume a traffic pattern  $b(f) = \sigma_i$ , where  $f$  is a flow that traverses a link in level  $i$ , but no link in level  $i - 1$ .<sup>6</sup> Assume that  $\sigma_1 = 1$  and  $\sigma_2 = \sigma_3 = \sigma$ .<sup>7</sup> Then, the traffic pattern is interference-free and the oversubscription parameter  $\omega(k, \sigma) = \frac{k}{2\beta(k, \sigma)}$  corresponds to the optimal design, where*

$$\beta(k, \sigma) = \left\lceil \frac{(k^4 - k^3)}{2(k^3 - k^2) + \sigma(2k^2 - 8)} \right\rceil \quad (10)$$

*is the number of deployed spine blocks and  $\lceil \cdot \rceil$  is the ceiling operator. Equivalently, a  $Clos(k, 3)$  with  $\beta$  spine blocks is optimal if  $\sigma_1 = 1$  and  $\sigma_2 = \sigma_3 = \sigma(k, \beta)$ , where:*

$$\sigma(k, \beta) = \begin{cases} \frac{(k^3 - k^2)(k - 2\beta)}{2\beta(k^2 - 4)}, & \text{for } 1 \leq \beta < k/2 \\ 1, & \text{for } \beta = k/2 \end{cases} \quad (11)$$

PROOF. See appendix A.8.  $\square$

Using Equation (10), Fig. 5 plots the value of the optimal number of spine blocks  $\beta$  to deploy in a  $Clos(k, 3)$  interconnect as a function of the inverse<sup>8</sup> of traffic skewness  $1/\sigma$  and for various values of the radix parameter  $k$ . For instance, for a folded-Clos with radix  $k = 48$  and traffic skewness  $\sigma = 10$  ( $\sigma^{-1} = 0.1$ ), it is enough to deploy  $\beta = 20$  spine blocks, each consisting of  $k/2 = 24$  switches. Since the  $Clos(48, 3)$  interconnect without oversubscription has a total of  $k/2 = 24$  spine blocks, this yields a total savings of 4 spine blocks or, equivalently, 96 spine switches (this optimal design is represented as a red dot in Fig. 5). As traffic becomes more skewed

<sup>6</sup>So a flow that goes through all levels of the folded-Clos transmits  $\sigma_1$ , and a flow between a pair of hosts that are nearest neighbors in the folded-Clos transmits  $\sigma_3$ .

<sup>7</sup>For simplicity, we assume all intrapod traffic is equally skewed, although it is also possible to derive the equations for the general case.

<sup>8</sup>Using the inverse of  $\sigma$  helps illustrate the asymptotic bounds in this plot.

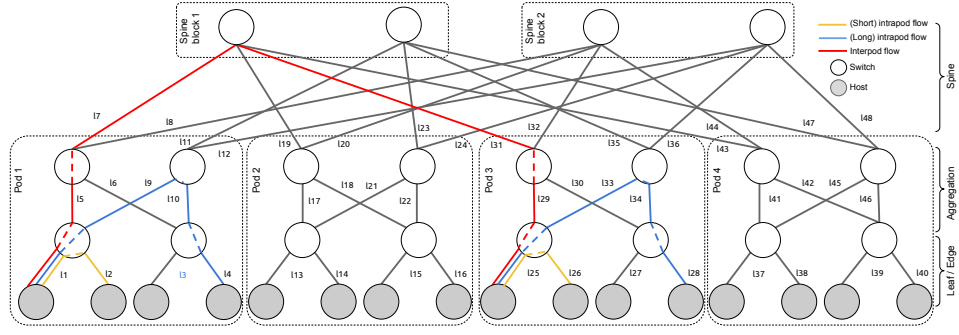


Figure 4: A folded-Clos network with 3 levels and radix 4, denoted as  $Clos(4, 3)$ .

( $\sigma^{-1} \rightarrow 0$ ) fewer spine blocks are needed. Interestingly, Fig. 5 also reveals that as  $k$  increases, the optimal skewness value has an asymptotic bound, represented with dashed vertical lines. The position of these asymptotic lines can be obtained from Equation (11) by taking the limit of  $k \rightarrow \infty$  and setting  $\beta = \frac{k}{2} - i$ , where  $i$  corresponds to the number of spine blocks that can be eliminated from the folded-Clos without incurring any performance penalty:

$$\lim_{k \rightarrow \infty} \sigma(k, \beta = \frac{k}{2} - i) = \lim_{k \rightarrow \infty} \frac{(k^3 - k^2)(k - 2(\frac{k}{2} - i))}{2(\frac{k}{2} - i)(k^2 - 4)} = 2i \quad (12)$$

The above equation provides a simple but succinct rule that all  $Clos(k, 3)$  interconnects should satisfy to avoid unnecessary investment costs in the provisioning of the spine layer, which we formalize as follows:

**COROLLARY 4.2.** *Minimum oversubscription requirement of  $Clos(k, 3)$ . If traffic skewness  $\sigma$  is larger than  $2(i + 1)$ , then  $i$  spine blocks can be removed from the interconnect without incurring any performance penalty, regardless of the radix value  $k$ . Equivalently, the interconnect can be oversubscribed by at least a factor  $\omega = \frac{k}{k - 2i}$  without incurring any performance penalty.*

Note that the above rule provides a loose bound since it is generally applicable to any radix value. For a known specific radix, designers can directly use Equation (10) to compute a more precise bound on the minimum oversubscription requirement. Next, we state the maximum flow completion time of a  $Clos(k, 3)$ :

**LEMMA 4.3.** *Network completion time of a 3-level folded-Clos with radix  $k$  and skewed traffic. Assume that every host sends  $\sigma$  bits of information to every other host located in the same pod and 1 bit of information to every other host located in a remote pod. The network completion time of a 3-level folded-Clos with radix  $k$  is:*

$$\mu(\beta, \sigma) = \begin{cases} \frac{k^2 + \sigma(k^3 - k^2) - 4}{\sigma(k^4 - k^3)}, & \text{if } \sigma \leq \sigma(k, \beta) \\ \frac{2c}{4\beta c}, & \text{otherwise} \end{cases} \quad (13)$$

where  $c$  is the capacity of each switch port.

**PROOF.** See appendix A.9. □

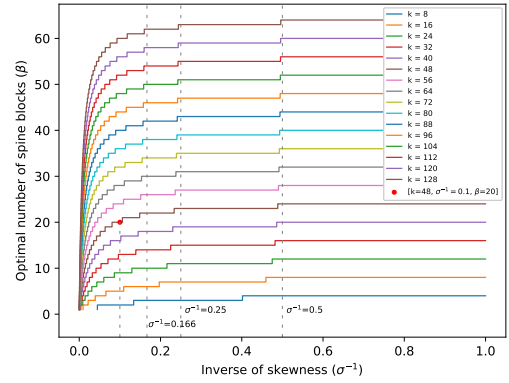


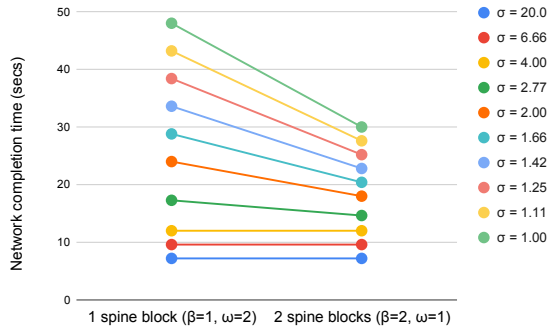
Figure 5: Optimal number of spine blocks as a function of traffic skewness for various radix values ( $k$ ).

Fig. 6 shows the network completion time<sup>9</sup> of  $Clos(4, 3)$  as a function of the number of spine blocks deployed  $\beta$  assuming a normalized link capacity of  $c = 1$  bps. (Note that from Equation (13), to obtain the maximum flow completion time for arbitrary values of  $c$ , we can simply scale the vertical axis by a factor of  $1/c$ .) This plot shows another interesting property of folded-Clos interconnects: as traffic skewness increases, the performance benefit of increasing the number of spine blocks diminishes, and no benefit at all is obtained after a certain threshold. In the case of  $Clos(4, 3)$ , this threshold is  $\sigma = 4$ , indicating that for all traffic patterns with skewness  $\sigma \geq 4$ , a network with one spine block is as performant as one with two spine blocks. Later in Section 5.2 we empirically validate this result.

In Appendix E we provide a plot of the network completion time as a function of the oversubscription parameter  $\omega$  for a production-scale folded-Clos [1] with radix  $k = 48$ .

<sup>9</sup>In this chart, network completion time has been normalized to (divided by)  $\sigma$  to better illustrate its asymptotic behavior, so the corresponding traffic pattern is  $b(f) = 1/\sigma$  for interpod flows and  $b(f) = 1$  for intrapod flows.





**Figure 6: In a  $Clos(4, 3)$  network with skewness  $\sigma \geq 4$ , using one spine block constitutes an optimal design.**

## 5 EXPERIMENTS

To experimentally demonstrate the accuracy of the QTBS model, we use *G2-Mininet* [19], a network emulation framework based on Mininet [22] with a set of extensions developed by our team to support the analysis of QTBS. Leveraging software defined networking (SDN), *G2-Mininet* enables the creation and analysis of topologies (such as fat-trees, folded-Clos and dragonflies, among others) using production TCP/IP code, including production-grade implementations of congestion control algorithms such as BBR, Cubic or Reno. (See Appendix H.) We are open sourcing *G2-Mininet* and all the experiments presented in this paper, hoping this will also enable the research community to verify our findings and further experiment with QTBS. More than 600 network simulations for a total of more than 800 hours of simulation time were used to verify the correctness of the model. We present a summary of these results in the following sections.

### 5.1 Experiments with Fat-Trees

In this first set of experiments our objective is to empirically demonstrate the existence of an optimal fat-tree design—a design that both minimizes network completion time and maximizes network throughput for a given cost. We start by simulating a  $FT(3, 2)$  interconnect—i.e., a fat-tree with two levels ( $L = 2$ ) and with each switch having three children  $n = 3$ . We connect every pair of nodes with two TCP flows (one for each direction), for a total of  $n^L(n^L - 1) = 72$  flows. In the first set of experiments, we assume uniform traffic ( $\sigma = 1$ ). Using Equation (4), we have that the optimal design corresponds to  $c_1 = n^2 \cdot c_2 / (n + 1) = 3^2 \cdot c_2 / (3 + 1) = 2.25 \cdot c_2$ .

Fig. 7 shows the result of simulating a variety of designs with  $c_2 = 5$  Mbps and  $c_1 \in \{2.5, 5, 7.5, 10, 11.25, 15, 20, 25\}$  Mbps—resulting in values for  $c_1/c_2$  of 0.5, 1, 1.5, 2, 2.25, 3, 4, 5. (Note that, without loss of generality, we could pick any value for  $c_2$  and scale  $c_1$  accordingly). Results are shown for the BBR, Cubic and Reno congestion control algorithms, and for both experimental and theoretical (according to QTBS)

values. The network completion time corresponds to the upper envelope of the curves—i.e., the maximum completion time of the interpod and intrapod flows for any value of  $c_1/c_2$ . As predicted by QTBS, the plots show that the optimal design is found at  $c_1/c_2 = 2.25$ . According to Theorem 2.10, this design wastes no bandwidth, minimizes network completion time and maximizes network throughput. For  $c_1/c_2 < 2.25$ , the completion time of the interpod flows increases while that of the intrapod flows decreases. Choosing a design in this region wastes bandwidth at the leaf links, increases network completion time (due to the longer completion time of the interpod flows) and decreases network throughput. As shown also, a design in the region  $c_1/c_2 > 2.25$  achieves the same network completion time and network throughput as the design  $c_1/c_2 = 2.25$ , regardless of how large the capacity of a spine link ( $c_1$ ) is. A design in this region wastes bandwidth at the spine links and is more costly than the optimal design, thus it should be avoided.

BBR almost perfectly follows the predictions by QTBS. Cubic and Reno also follow them but not as accurately. This is also reflected in Fig. 8 where Jain’s fairness index is presented [15]. As shown in [27], this index captures how accurately an experiment is able to match the behavior according to the QTBS model. The index ranges from 0 to 1, with values close to 1 indicating high accuracy. In Fig. 8, Jain’s index for BBR is close to 1, while Cubic and Reno follow the model slightly less accurately—although they still qualitatively behave according to the model. This finding is aligned with the results in [27], and is explained by the advanced congestion-based techniques used by BBR [7], which allow it to more accurately infer the optimal transmission rate of each flow.

Fig. 9 shows the cumulative distribution function of all the experiments run in Fig. 7a. The figure shows how increasing  $c_1$  helps reduce the maximum flow completion time until the optimal value  $c_1 = 2.25 \cdot c_2$  is reached (green line). At this point, all completion times are equalized (Theorem 2.10) and increasing  $c_1$  beyond this value does not qualitatively alter the completion time of the flows.

In Fig. 11 we present experiments with a skewed traffic pattern consisting of  $b(f) = 1$  for interpod flows and  $b(f) = \sigma$  for intrapod flows, where  $\sigma \in \{2, 4, 10\}$ . Using Equation (8), we have that the optimal design satisfies  $c_1 = n^2 \cdot c_2 / (n + \sigma) = 9 \cdot c_2 / (3 + \sigma)$ . This leads to three optimal designs, one for each traffic pattern:  $c_1 = 1.8 \cdot c_2$  for  $\sigma = 2$ ,  $c_1 = 1.2857 \cdot c_2$  for  $\sigma = 4$ , and  $c_1 = 0.69231 \cdot c_2$  for  $\sigma = 10$ . As shown in the figure, the green and the red lines—which correspond to the experimental average completion time of the interpod and intrapod flows—cross each other right at the optimal design point, marked with a vertical dashed line. As it was shown in Fig. 3, the more skewed traffic is, the higher we can reduce the capacity of the spine links. This is shown in Fig. 11 with the crossing points of the completion time curves for

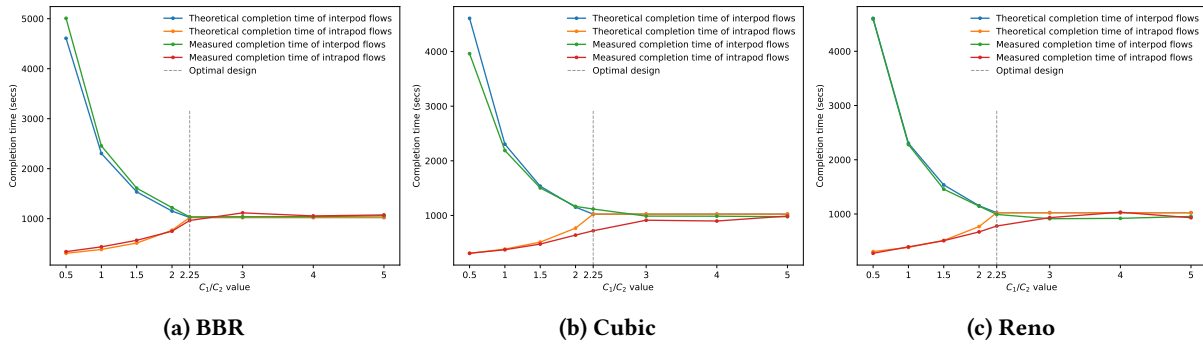


Figure 7: Flow completion time for  $FT(3, 2)$  and uniform traffic.

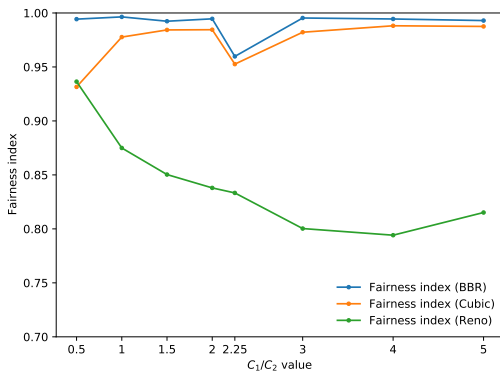


Figure 8: Jain's fairness index for  $FT(3, 2)$ .

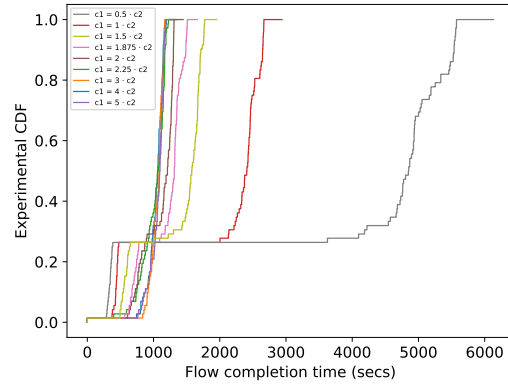


Figure 9: CDFs of flow completion time for  $FT(3, 2)$ .

interpod and intrapod flows shifting to the left as  $\sigma$  increases. Once again the optimal design does not waste any bandwidth (Theorem 2.10), since all flows terminate at the same time and for the whole duration of the simulation every link is saturated. A difference with the uniform case in Fig. 7 is that, for  $c_1 > 9 \cdot c_2 / (3 + \sigma)$ , interpod flows continue to reduce their completion time. However, intrapod flows' completion time stays flat, thus the network completion time stays flat too. An interesting *ripple effect* is produced in the bottleneck structure that leads to this behavior: as  $c_1$  increases, interpod flows receive more bandwidth, but in doing so they need to steal some of it from the intrapod flows, since they both share the intrapod links. At the same time, the interpod flows will finish sooner as they get more bandwidth, and so they will free bandwidth for the intrapod flows sooner, creating this time a positive impact for them. The two effects (one negative, one positive) exactly cancel out, and thus intrapod flows don't see their completion time affected as  $c_1$  increases beyond the optimal point.

### 5.2 Experiments with Folded-Clos

In this section we focus on empirically demonstrating the behavior of an optimal folded-Clos design according to the

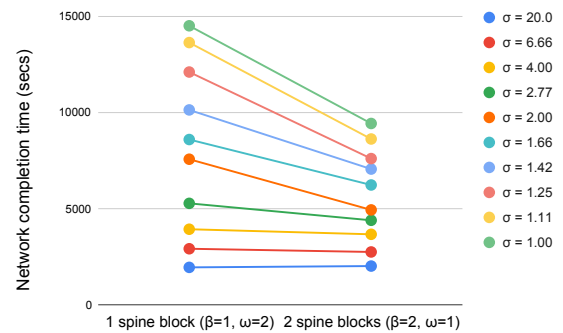


Figure 10: Experimental demonstration that an oversubscription of 2:1 is an optimal design for  $Clos(4, 3)$  and  $\sigma \geq 4$ .

equations presented in Section 4. We use *G2-Mininet* to emulate a  $Clos(4, 3)$  (see Fig. 4), consisting of 20 switches of radix 4, 48 links, 4 pods and 16 hosts. Every pair of hosts is connected with two flows (one for each direction), for a total

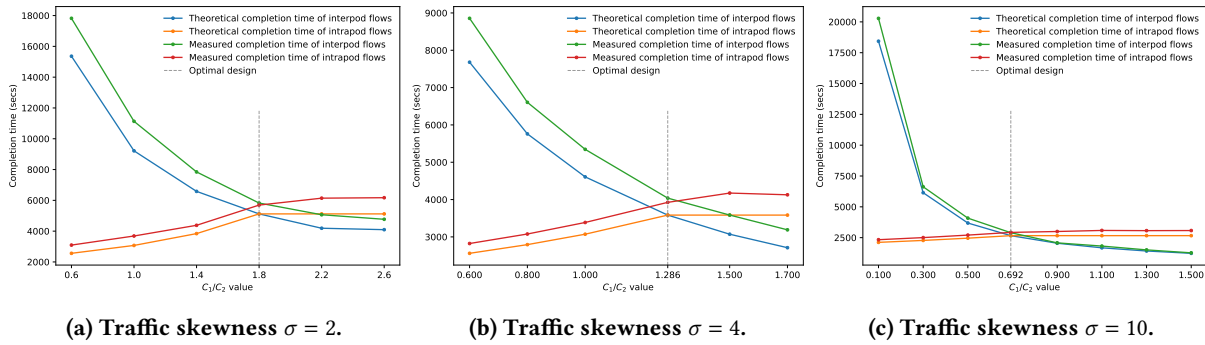


Figure 11: Flow completion time for  $FT(3, 2)$  using BBR.

of 240 flows. Fig. 10 presents the results of running experiments for the two possible  $Clos(4, 3)$  designs—with  $\beta = 1$  and  $\beta = 2$  spine blocks—and for a skewed traffic pattern consisting of  $b(f) = \sigma_1$  for interpod flows and  $b(f) = \sigma_2$  for intrapod flows, where  $\sigma_2 = 64$  MB,  $\sigma_1 = \sigma_2/\sigma$  Mbps and  $\sigma \in \{1, 1.11, 1.25, 1.42, 1.66, 2, 2.77, 4, 6.66, 20\}$ . All experiments in this chart were performed using BBR. As shown, the interconnect behaves as predicted by QTBS in Fig. 6. Using a 1-spine block configuration (thus reducing the cost of the interconnect’s spine level by half) is an optimal design for  $\sigma \geq 4$ , since using the 2-spine block configuration does not improve the maximum completion time and wastes bandwidth at the spine level. For  $\sigma < 4$ , both designs waste bandwidth, so the choice for one design or the other has to be based on the budget constraints and the performance objectives of the application.

In Fig. 12 we compare the network completion time obtained from the experiments against the value projected by QTBS according to Equation (13) and for the case  $\beta = 1$  (a  $\omega=2:1$  oversubscribed configuration). The plot shows that the model is able to describe the actual behavior of the interconnect fairly accurately. The model always lays below the experimental values due to imperfections of the congestion control algorithm that regulates the transmission rate of each flow, which leads to slightly higher network completion time than the theoretical one. (This behavior of the QTBS model is also explained in [27].) We present additional experiments supporting the validation of the QTBS model in Appendix I.

## 6 RELATED WORK

Data center networks have been the subject of intense research in the networking community. While many different topologies have been proposed and studied, in our work we focus on three of the most widely used interconnects: fat-trees, folded-Clos and dragonflies. Leiserson [20] demonstrated that fat-trees are universally efficient interconnects. His work, however, did not take into account the effects of the congestion-control algorithms that are part of modern

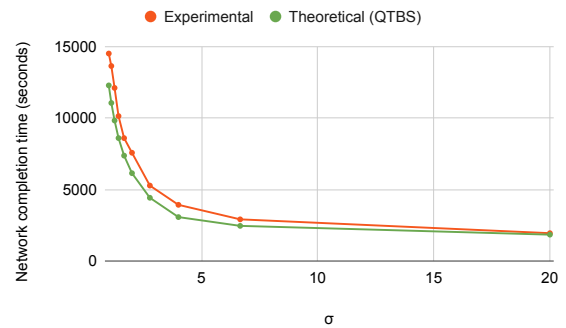


Figure 12: Comparison of experimental network completion time with QTBS model for  $\beta = 1$ .

communication networks. In this paper, we show that, for congestion-controlled networks, full fat-trees are only cost-efficient when more traffic is sent between pairs of hosts that are farther apart than between hosts that are closer together. This implies that full fat-trees should not be used in interconnects that transport traffic with locality patterns. We also contribute to the understanding of fat-trees performance by identifying the set of link capacities that make a fat-tree optimal for a given traffic pattern.

Folded-Clos is the dominant topology in large-scale data centers. The influential work by Al-Fares et al. [1] demonstrated the scalability and cost-effectiveness of this class of interconnects. This and follow up work has produced an extensive literature around the subject of design and capacity planning driven by empirical production-scale experiments (e.g., [16, 28, 30]). To the best of our knowledge, however, our work is the first to provide a formal model that can help network architects identify optimal oversubscription configurations of a folded-Clos for a given traffic pattern.

In [31], Singla et al. provide an upper bound on network throughput for homogeneous topologies with the assumption that all switches are identical, and then experimentally

demonstrate that random graphs achieve throughput close to the bound. Our work differs and complements this work in two ways. First, we provide a general mathematical model of data centers that is applicable to both homogeneous and heterogeneous interconnects and reveals the designs that are simultaneously optimal in throughput and latency. Secondly, while in this work we focus on deterministic, well-structured topologies that allow us to derive closed-form symbolic equations of the optimal designs, it is also possible to use QTBS to perform non-symbolic numerical analysis [26]. Left as future work, this opens up the possibility of using the proposed framework to numerically model and study the performance of unstructured topologies such as random graphs.

## 7 ASSUMPTIONS AND GENERALIZATIONS

In our model, an optimal design is one that outperforms (in network completion time and throughput) any other design with the same or less cost, where cost corresponds to total capacity. That is, it is impossible to outperform the optimal design by shifting capacity from one region of the network to another, or by reducing the total aggregated link capacity. QTBS however is general, and it could also be applied to identify optimal designs based on other definitions of cost.

Our model also assumes that routing is fixed. We see our work as a first step towards enabling future generalizations that accommodate for dynamic routing. In production networks, the global optimal routing solution can change every time a new flow arrives or departs from the network. It is interesting to note that for interference-free traffic patterns, the optimal design has all flows terminating at the same time. Thus, optimal routing remains static in this case throughout the full execution of a batch. We leave for future work the problem of extending the QTBS model to account for dynamic routing of non-interference free traffic patterns.

We currently assume that the traffic pattern is stable. In production networks, traffic instead consists of the superposition of many different traffic patterns. In this case, network operators may approximate their problem by a representative or average traffic pattern to take advantage of the power of our formal model. In addition, while the analysis presented in this paper is purely symbolic, in future work QTBS could also be used to develop numerical methods and broaden the applicability of our work to explicitly take into account variation in the traffic pattern, much as physicists do to make calculations when no closed-form solution exists. Finally, we can use our framework to design optimal overlay networks in addition to the underlying hardware network. This would allow us to handle multiple jobs with different traffic patterns by assigning each to a different overlay and to address changes in the workload of the network by reconfiguring these overlays. This opens up the possibility of using QTBS to address the multi-tenancy bandwidth allocation and

network slicing problems (see for instance [3] and [12]). We leave these interesting direction for future work.

We show that the assumption of interference-freeness is reasonable in that it corresponds to the important set of patterns that display data locality properties. However, we note that even if the traffic pattern is not interference free, QTBS can still be used as a model to design interconnects. Similar to the case of dynamic traffic patterns, in this case operators could use QTBS to perform non-symbolical numerical analysis to compare the various designs without having to resort to costly trial-and-error implementations.

While our work focuses on optimizing performance, the QTBS model is generic and, in future work, it can also be used to address other important design criteria. Consider for instance *resiliency analysis*. Using the predictive capabilities of QTBS, an operator can use the proposed model to measure the negative effect of a link failure, and provision additional redundant capacity or links to protect those regions of the network whose failure would lead to higher impact.

## 8 CONCLUSIONS

We present a model to compute optimal link capacity settings in data center networks based on the recently introduced Quantitative Theory of Bottleneck Structures (QTBS). Using this framework, we show that for interconnects with traffic patterns that are *interference-free*, there exists a design that optimizes both the completion time and throughput needed to execute such traffic patterns. We demonstrate that fat-trees, folded-Clos and dragonfly topologies satisfy this property for typical production traffic patterns that display data locality properties. In particular, we identify the set of designs that are always inefficient (more costly without yielding any performance benefit), regardless of the traffic pattern. QTBS proves that if the traffic pattern is such that more data is transmitted between two hosts in the same pod than two hosts in different pods (which holds in most production interconnects), proportional designs are optimal in both latency and throughput. Moreover, we demonstrate that increasing the capacity of the interpod links (the most expensive layer in modern data center interconnects) above the value determined by this optimal design yields no performance gain. These results provide an engineering framework to help network architects and operators identify the right size of an interconnect that delivers maximum performance while avoiding unnecessary costs. Forthcoming work will focus on applying QTBS to production-level traffic patterns and validating these results in production interconnects.

### Acknowledgments

We wish to thank the shepherd Radhika Niranjana Mysore for her insightful comments and guidance as well as the anonymous reviewers for their valuable feedback. This work was partly supported by the DOE DE-SC0019523.

## REFERENCES

- [1] Mohammad Al-Fares, Alexander Loukissas, and Amin Vahdat. 2008. A Scalable, Commodity Data Center Network Architecture. In *Proceedings of the ACM SIGCOMM 2008 Conference on Data Communication (SIGCOMM '08)*. Association for Computing Machinery, New York, NY, USA, Article 1, 12 pages. <https://doi.org/10.1145/1402958.1402967>
- [2] Mohammad Alizadeh, Tom Edsall, Sarang Dharmapurikar, Ramanan Vaidyanathan, Kevin Chu, Andy Fingerhut, Vinh The Lam, Francis Matus, Rong Pan, Navindra Yadav, and George Varghese. 2014. CONGA: Distributed Congestion-Aware Load Balancing for Datacenters. In *Proceedings of the 2014 ACM Conference on SIGCOMM (SIGCOMM '14)*. Association for Computing Machinery, New York, NY, USA, Article 1, 12 pages. <https://doi.org/10.1145/2619239.2626316>
- [3] Hitesh Ballani, Paolo Costa, Thomas Karagiannis, and Ant Rowstron. 2011. Towards predictable datacenter networks. In *Proceedings of the ACM SIGCOMM 2011 Conference*. 242–253.
- [4] Theophilus Benson, Aditya Akella, and David A. Maltz. 2010. Network Traffic Characteristics of Data Centers in the Wild. In *Proceedings of the 10th ACM SIGCOMM Conference on Internet Measurement (IMC '10)*. Association for Computing Machinery, New York, NY, USA, Article 1, 14 pages. <https://doi.org/10.1145/1879141.1879175>
- [5] Theophilus Benson, Ashok Anand, Aditya Akella, and Ming Zhang. 2010. Understanding data center traffic characteristics. *ACM SIGCOMM Computer Communication Review* 40, 1 (2010), 92–99.
- [6] Dimitri P. Bertsekas and Robert G. Gallager. 1992. *Data Networks*. Vol. 2. Prentice-Hall Inc., Englewood Cliffs, New Jersey 07632.
- [7] Neal Cardwell, Yuchung Cheng, C. Stephen Gunn, Soheil Hassas Yeganeh, and Van Jacobson. 2016. BBR: Congestion-Based Congestion Control. *ACM Queue* 14, 5, Article 50 (October 2016), 34 pages. <https://doi.org/10.1145/3012426.3022184>
- [8] Dah-Ming Chiu and Raj Jain. 1989. Analysis of the increase and decrease algorithms for congestion avoidance in computer networks. *Computer Networks and ISDN systems* 17, 1 (1989), 1–14.
- [9] Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. 2009. *Introduction to Algorithms, Third Edition* (3rd ed.). The MIT Press, Cambridge, Massachusetts.
- [10] Wolfgang E Denzel, Jian Li, Peter Walker, and Yuho Jin. 2010. A framework for end-to-end simulation of high-performance computing systems. *Simulation* 86, 5-6 (2010), 331–350.
- [11] Nandita Dukkipati and Nick McKeown. 2006. Why flow-completion time is the right metric for congestion control. *ACM SIGCOMM Computer Communication Review* 36, 1 (2006), 59–62.
- [12] Andreas Fischer, Juan Felipe Botero, Michael Till Beck, Hermann De Meer, and Xavier Hesselbach. 2013. Virtual network embedding: A survey. *IEEE Communications Surveys & Tutorials* 15, 4 (2013), 1888–1906.
- [13] Iperf.fr. 2019. iPerf - The ultimate speed test tool for TCP, UDP and SCTP. (2019). Retrieved October 24, 2019 from <https://iperf.fr/>
- [14] Van Jacobson. 1988. Congestion Avoidance and Control. *SIGCOMM computer communication review* 18, 4 (August 1988), 314–329. <https://doi.org/10.1145/52325.52356>
- [15] Raj Jain, Dah-Ming W. Chiu, and William R. Hawe. 1998. A Quantitative Measure Of Fairness And Discrimination For Resource Allocation In Shared Computer Systems. *CoRR cs.NI/9809099* (1998), 38. <http://arxiv.org/abs/cs.NI/9809099>
- [16] Sangeetha Abdu Jyothi, Ankit Singla, P Brighten Godfrey, and Alexandra Kolla. 2016. Measuring and understanding throughput of network topologies. In *SC'16: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*. IEEE, 761–772.
- [17] Srikanth Kandula, Sudipta Sengupta, Albert Greenberg, Parveen Patel, and Ronnie Chaiken. 2009. The nature of data center traffic: measurements & analysis. In *Proceedings of the 9th ACM SIGCOMM conference on Internet measurement*. Association for Computing Machinery, New York, NY, USA, 202–208.
- [18] John Kim, William J Dally, Steve Scott, and Dennis Abts. 2008. Technology-driven, highly-scalable dragonfly topology. In *2008 International Symposium on Computer Architecture*. IEEE, 77–88.
- [19] Reservoir Labs. 2021. G2-Mininet: Mininet extensions to support the analysis of the bottleneck structure of networks. <https://github.com/reservoirlabs/g2-mininet>. (2021). Retrieved June 11, 2021 from <https://github.com/reservoirlabs/g2-mininet>
- [20] Charles E Leiserson. 1985. Fat-trees: universal networks for hardware-efficient supercomputing. *IEEE transactions on Computers* 100, 10 (1985), 892–901.
- [21] George Michelogiannakis, Yiwen Shen, Min Yee Teh, Xiang Meng, Benjamin Aivazi, Taylor Groves, John Shalf, Madeleine Glick, Manya Ghobadi, Larry Dennison, and Keren Bergman. 2019. Bandwidth Steering in HPC Using Silicon Nanophotonics. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC '19)*. Association for Computing Machinery, New York, NY, USA, Article Article 41, 25 pages. <https://doi.org/10.1145/3295500.3356145>
- [22] Mininet. 2019. Mininet: An Instant Virtual Network on your Laptop (or other PC). (2019). Retrieved October 24, 2019 from <http://mininet.org/>
- [23] A. Munir, I. A. Qazi, Z. A. Uzmi, A. Mushtaq, S. N. Ismail, M. S. Iqbal, and B. Khan. 2013. Minimizing flow completion times in data centers. In *2013 Proceedings IEEE INFOCOM*. IEEE, New York, NY, USA, 2157–2165. <https://doi.org/10.1109/INFCOM.2013.6567018>
- [24] Peter Phaal, Sonia Panchen, and Neil McKee. 2001. sFlow Specifications. InMon Corporation. *IETF RFC 3176* (2001).
- [25] NOX Repo POX. 2019. The POX Network Software Platform. (2019). <https://noxrepo.github.io/pox-doc/html/>, Last accessed 09/24/2019.
- [26] Jordi Ros-Giralt, Noah Amsel, Sruthi Yellamraju, James Ezick, Richard Lethin, Yuang Jiang, Aosong Feng, Leandros Tassiulas, Zhenguo Wu, Min Yeh Teh, and Keren Bergman. 2021. A Quantitative Theory of Bottleneck Structures for Data Networks. *IEEE Transactions on Networking (under review)* (2021).
- [27] Jordi Ros-Giralt, Atul Bohara, Sruthi Yellamraju, M. Harper Langston, Richard Lethin, Yuang Jiang, Leandros Tassiulas, Josie Li, Yuanlong Tan, and Malathi Veeraraghavan. 2019. On the Bottleneck Structure of Congestion-Controlled Networks. *Proc. ACM Meas. Anal. Comput. Syst.* 3, 3, Article Article 59 (Dec. 2019), 31 pages. <https://doi.org/10.1145/3366707>
- [28] Arjun Roy, Hongyi Zeng, Jasmeet Bagga, George Porter, and Alex C Snoeren. 2015. Inside the social network's (datacenter) network. In *Proceedings of the 2015 ACM Conference on Special Interest Group on Data Communication*. 123–137.
- [29] Daniele Sensi, Salvatore Girolamo, Kim McMahon, Duncan Roweth, and Torsten Hoefer. 2020. An In-Depth Analysis of the Slingshot Interconnect. In *2020 SC20: International Conference for High Performance Computing, Networking, Storage and Analysis (SC)*. IEEE Computer Society, 481–494.
- [30] Arjun Singh, Joon Ong, Amit Agarwal, Glen Anderson, Ashby Armistead, Roy Bannon, Seb Boving, Gaurav Desai, Bob Felderman, Paulie Germano, Anand Kanagala, Jeff Provost, Jason Simmons, Eiichi Tanda, Jim Wanderer, Urs Hölzle, Stephen Stuart, and Amin Vahdat. 2015. Jupiter Rising: A Decade of Clos Topologies and Centralized Control in Google's Datacenter Network. In *Proceedings of the 2015 ACM Conference on Special Interest Group on Data Communication (SIGCOMM '15)*. Association for Computing Machinery, New York, NY, USA, 183–197. <https://doi.org/10.1145/2785956.2787508>



- [31] Ankit Singla, P Brighten Godfrey, and Alexandra Kolla. 2014. High throughput data center topology design. In *11th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 14)*. 29–41.
- [32] David Zats, Tathagata Das, Prashanth Mohan, Dhruva Borthakur, and Randy Katz. 2012. DeTail: Reducing the Flow Completion Time Tail in Datacenter Networks. In *Proceedings of the ACM SIGCOMM 2012 Conference on Applications, Technologies, Architectures, and Protocols for Computer Communication (SIGCOMM '12)*. Association for Computing Machinery, New York, NY, USA, Article 1, 12 pages. <https://doi.org/10.1145/2342356.2342390>

Appendices are supporting material that has not been peer-reviewed.

## 9 APPENDICES

### A MATHEMATICAL PROOFS

#### A.1 Proof of Theorem 2.6

**Optimality of wasteless designs.** *For a given topology and traffic pattern, if a design  $c$  is wasteless, then it is impossible to improve on the network completion time or network throughput of the design without adding more capacity. That is, if  $c'$  is an alternative design for which  $\mu(b, c') < \mu(b, c)$  or  $T(b, c') > T(b, c)$ , then  $\sum_{l \in \mathcal{L}} c'(l) > \sum_{l \in \mathcal{L}} c(l)$ .*

**PROOF.** We first prove the statement regarding network completion time. The total amount of data transmitted by any flow  $f$  over the course of the transmission is  $b(f)$ . So if  $r_c(f, t)$  is the rate of flow  $f$  at time  $t$  under design  $c$ , then

$$b(f) = \int_0^{\mu(b,c)} r_c(f, t) dt$$

Taking sums on both sides:

$$\sum_{l \in \mathcal{L}} \sum_{f \in \mathcal{F}_l} b(f) = \int_0^{\mu(b,c)} \sum_{l \in \mathcal{L}} \sum_{f \in \mathcal{F}_l} r_c(f, t) dt$$

where  $\mathcal{F}_l$  is the set of flows that use link  $l$ . The same holds for the alternative design  $c'$ :

$$\sum_{l \in \mathcal{L}} \sum_{f \in \mathcal{F}_l} b(f) = \int_0^{\mu(b,c')} \sum_{l \in \mathcal{L}} \sum_{f \in \mathcal{F}_l} r_{c'}(f, t) dt$$

As always, the rates must respect the capacity constraint of each link:

$$\sum_{f \in \mathcal{F}_l} r_{c'}(f, t) \leq c'(l)$$

Combining this with the equation above,

$$\sum_{l \in \mathcal{L}} \sum_{f \in \mathcal{F}_l} b(f) \leq \int_0^{\mu(b,c')} \sum_{l \in \mathcal{L}} c'(l) dt = \mu(b, c') \cdot \sum_{l \in \mathcal{L}} c'(l)$$

By assumption, the original design  $c$  is wasteless, so the full capacity of each link is used:

$$\sum_{f \in \mathcal{F}_l} r_c(f, t) = c(l)$$

Thus for  $c$ , we derive an equality instead of an inequality:

$$\sum_{l \in \mathcal{L}} \sum_{f \in \mathcal{F}_l} b(f) = \mu(b, c) \cdot \sum_{l \in \mathcal{L}} c(l)$$

Combining this with the inequality involving  $c'$  above,

$$\mu(b, c) \cdot \sum_{l \in \mathcal{L}} c(l) \leq \mu(b, c') \cdot \sum_{l \in \mathcal{L}} c'(l)$$

Thus, if  $\mu(b, c') < \mu(b, c)$ , then  $\sum_{l \in \mathcal{L}} c(l) < \sum_{l \in \mathcal{L}} c'(l)$ . If no bandwidth is wasted for one batch, the scheduling  $n$  batches sequentially wastes no bandwidth either. Thus the same argument proves the statement about network throughput.  $\square$

#### A.2 Proof of Theorem 2.8

**Non-proportional designs waste bandwidth.** *If a design wastes no bandwidth on a given traffic pattern, then it is the proportional design for the traffic pattern.*

**PROOF.** Let  $c$  be the design, and let  $b$  be the traffic pattern. Let  $\mu$  be the completion time of the network. The design wastes no bandwidth, so for all links  $l$  and all  $t \in [0, \mu]$ ,

$$\sum_{f \in \mathcal{F}_l} r_c(f, t) = c(l)$$

where  $\mathcal{F}_l$  is the set of flows that traverse link  $l$ . The total amount of data transmitted by a flow  $f$  during  $[0, \mu]$  is exactly  $b(f)$  (i.e., all of its data), so taking the integral of the left side of the above equation,

$$\int_0^\mu \sum_{f \in \mathcal{F}_l} r_c(f, t) dt = \sum_{f \in \mathcal{F}_l} \int_0^\mu r_c(f, t) dt = \sum_{f \in \mathcal{F}_l} b(f)$$

Integrating the right side of the equality above,

$$\int_0^\mu c(l) dt = \mu \cdot c(l)$$

These two expressions are equal. Thus,

$$c(l) = \frac{1}{\mu} \sum_{f \in \mathcal{F}_l} b(f)$$

and  $c$  is proportional.  $\square$

#### A.3 Proof of Theorem 2.10

We first prove the following lemma:

**LEMMA A.1.** *If the traffic pattern is interference free, all flows start transmitting at the same time, and the proportional design is used, then each flow's rate is proportional to its size:  $r_f = \alpha \cdot b(f)$ , where  $\alpha$  is the coefficient of the proportional design.*

PROOF. We prove the lemma by induction on the flows  $f$ , going in order from smallest to largest (if there are repeated sizes, this order may not be unique but the argument still holds). The induction hypothesis is that all flows that precede  $f$  in this ordering have rates proportional to their size. This hypothesis holds for the first flow, since there are no preceding flows. Now let  $f$  be any flow. We must prove that if the induction hypothesis holds, then  $r_f = \alpha \cdot b(f)$ . Since  $b$  is interference-free, there exists a link  $l$  such that  $f$  traverses  $l$ , and for all flows  $f'$  that traverse  $l$ ,  $b(f') \leq b(f)$ . Since  $c$  is proportional,

$$c(l) = \alpha \sum_{f' \in \mathcal{F}_l} b(f')$$

By the induction hypothesis, all flows for which  $b(f') < b(f)$  have rates proportional to their size. Thus, the remaining capacity of  $l$  is

$$c(l) - \sum_{\substack{f' \in \mathcal{F}_l \\ b(f') < b(f)}} \alpha \cdot b(f') = \alpha \sum_{\substack{f' \in \mathcal{F}_l \\ b(f') = b(f)}} b(f')$$

This capacity is divided among the remaining flows that traverse  $l$ , which all have the same size as  $f$ . Thus, the fairshare of link  $l$  is  $\alpha \cdot b(f)$ , and the transmission rate of flow  $f$  is  $\alpha \cdot b(f)$  too. By induction, the rate of any flow  $f$  is  $r_f = \alpha \cdot b(f)$ .  $\square$

We now prove Theorem 2.10:

**Interference-free proportional designs are wasteful.** *If the traffic pattern is interference free and all flows start transmitting at the same time, then the proportional design wastes no bandwidth, and furthermore, all flows finish transmitting at the same time.*

PROOF. By Lemma A.1, for all  $f \in \mathcal{F}$ ,  $r_f = \alpha \cdot b(f)$  where  $\alpha$  is the coefficient of the proportional design. For any link  $l$ , the unused capacity at the beginning of the transmission is

$$c(l) - \sum_{f \in \mathcal{F}_l} r_f = c(l) - \sum_{f \in \mathcal{F}_l} \alpha \cdot b(f) = 0$$

Thus, at the beginning of the transmission, the capacity of each link is fully utilized. Since each flow's rate is proportional to its size, all flows finish simultaneously. Thus, none of the flows' rates change during the transmission. Thus, throughout the transmission, the capacity of each link is fully utilized.  $\square$

#### A.4 Proportional Designs Yield Maximal Throughput

**THEOREM A.2.** *For a given topology and traffic pattern, if the design of the network is proportional, then it is impossible to improve network throughput without adding more capacity to the network.*

PROOF. We will show that if the design is proportional, there exists some integer  $B$  and some schedule for  $B$  batches such that no bandwidth is wasted. By Theorem 2.6, this shows that any alternative design with a better completion time for the  $n$  batches must have more capacity than the proportional design, and continues to hold as the number of batches grows without bound. Assume without loss of generality that the size of each flow  $b(f)$  is an integer. (In this discussion, we use the word “flow” to mean one of the flows of a *single* batch, not one of multiple copies of these flows from different batches.) Let  $B$  be the least common multiple of  $\{b(f) \mid f \in \mathcal{F}\}$ . If  $\alpha$  is the scale parameter of the proportional design, then we will transmit  $B$  batches according to the following schedule:

*For each flow  $f$ , begin transmitting  $b(f)$  copies of  $f$  every  $b(f)/\alpha$  seconds until time  $B/\alpha$ .*

We now prove that for a proportional design, this schedule wastes no bandwidth. Since  $c$  is proportional, for each link  $l$ , the capacity is

$$c(l) = \alpha \cdot \sum_{f \in \mathcal{F}_l} b(f)$$

For each  $f \in \mathcal{F}_l$ , we begin transmitting  $b(f)$  copies at time 0. Thus, the initial fairshare of link  $l$  is

$$\frac{c(l)}{\sum_{f \in \mathcal{F}_l} b(f)} = \alpha$$

This is true for all links. Thus the initial transmission rate of each flow is  $\alpha$ . This means that no bandwidth is wasted; since there are  $b(f)$  copies of each flow traversing each link,

$$\sum_{f \in \mathcal{F}_l} b(f) r_c(f, t) = \sum_{f \in \mathcal{F}_l} b(f) \cdot \alpha = c(l)$$

Thus, all  $b(f)$  copies of flow  $f$  will finish transmitting at time  $b(f)/\alpha$ . According to our schedule, these are immediately replaced by  $b(f)$  more copies of  $f$ . Thus, for as long as the network transmits, the rate of each flow is  $\alpha$ , and no bandwidth is wasted. At time  $B/\alpha$ , the number of times that we will have begun new transmissions of flow  $f$  is  $B/\alpha \div b(f)/\alpha = B/b(f)$  (this is an integer since  $b(f)$  divides  $B$  by definition). During each of these periods, we transmit  $b(f)$  copies of  $f$ , so by the end we will have finished transmitting  $B$  copies. Since we have completed  $B$  copies of every flow, this schedule successfully transmits  $B$  batches in the time  $B/\alpha$  without wasting any bandwidth. Furthermore, the throughput (completed batches per second) is  $B \div B/\alpha = \alpha$ .  $\square$

#### A.5 Proof of Lemma 3.1

**Optimal fat-tree with uniform traffic.** *Consider a generalized fat-tree  $FT([n_1, n_2, \dots, n_L])$  and let  $c_i$  be the capacity of its links at level  $i$ , for  $1 \leq i \leq L$ . Then, if the traffic pattern is uniform, it is interference-free and the following design is optimal:*

$$c_i = \frac{\rho_i}{\rho_L} c_L \quad (14)$$

where

$$\rho_i = 2 \left( \prod_{j=1}^i n_j - 1 \right) \prod_{j=i+1}^L n_j^2 \quad (15)$$

**PROOF.** We first derive an expression for  $\rho_i$ , the number of flows traversing each link in level  $i$ . Each edge divides the hosts into two groups—the set of descendants of the edge, and the complement of that set. Let  $d_i$  be the number of hosts that are descendants of each link in level  $i$  of the tree. Then

$$d_i = \prod_{j=i+1}^L n_j$$

We now count the number of flows that traverse each link in level  $i$ . For a given link, two flows traverse it for each pair of one host that is descendant of the link and one host that is not a descendant (two flows because there are two possible directions along that path). The total number of hosts in the tree is

$$\prod_{j=1}^L n_j$$

Thus

$$\begin{aligned} \rho_i &= 2 \left( \prod_{j=1}^L n_j - d_i \right) d_i = 2 \left( \frac{1}{d_i} \prod_{j=1}^L n_j - 1 \right) d_i^2 = \\ &= 2 \left( \prod_{j=1}^i n_j - 1 \right) \prod_{j=i+1}^L n_j^2 \end{aligned}$$

We now show that the design given in the lemma is optimal. First, note that the given traffic pattern is trivially interference-free, since all flows are the same size. Next, note that the design is proportional. Each link in level  $i$  has capacity

$$c_i = \frac{\rho_i}{\rho_L} c_L$$

and is traversed by  $\rho_i$  flows of equal size. Thus, the capacity of each link is proportional to the sum of the sizes of the flows traversing it. Thus, by Theorem 2.10 the design is optimal with respect to completion time and throughput.  $\square$

## A.6 Proof of Lemma 3.2

**Optimal fat-tree with skewed traffic.** Consider a generalized fat-tree  $FT([n_1, n_2, \dots, n_L])$  and let  $c_i$  be the capacity of its links at level  $i$ , for  $1 \leq i \leq L$ . Assume a traffic pattern  $b(f) = \sigma_i$ , where  $f$  is a flow that traverses a link in level  $i$ , but no link in level  $i-1$ . If  $i \leq j \implies \sigma_i \leq \sigma_j$  for all  $1 \leq i, j \leq L$ , then the traffic pattern is interference free and

the following design is optimal:

$$c_1 = \pi_1 \sigma_1 \quad (16)$$

$$c_i = \pi_i \sigma_i + \frac{c_{i-1}}{n_i} \quad (17)$$

where

$$\pi_i = (n_i - 1) \prod_{j=i+1}^L n_j^2 \quad (18)$$

**PROOF.** We begin by deriving the expression for  $\pi_i$ , which is the number of flows that traverse each link in level  $i$  but do not traverse any link in level  $i-1$ . According to our definition of the traffic pattern, these flows will all be of size  $\sigma_i$ . The number of hosts that are descendants of each link in level  $i$  is

$$d_i = \prod_{j=i+1}^L n_j$$

(See the proof of Lemma 3.1.) Any given link in level  $i$  has a single parent link in level  $i-1$ . For a flow to be counted in  $\pi_i$ , it must traverse the given link, but not its parent. That is, it must have one endpoint that is a descendant of the given link, and one endpoint that is a descendant of the given link's siblings. Since the given link is in level  $i$ , it has  $i$  siblings, and each of them has  $d_i$  descendants. Thus,

$$\pi_i = 2d_i \cdot (n-1)d_i = 2(n-1) \prod_{j=i+1}^L n_j^2$$

(We include the factor of two because there are two possible flows for each pair of hosts, one in each direction).

We now show that the specified design is optimal. Note that while multiple designs specify the given equations, they are all equivalent up to scaling.

First note that because of the assumption  $\sigma_i \leq \sigma_j$ , the traffic pattern is interference free. The flows of size  $\sigma_i$  traverse links of level  $i$ , which are only traversed by other flows of smaller or equal size  $\sigma_1, \dots, \sigma_i$ .

Second, we prove by induction that the given design is proportional. By definition, the number of flows that traverse each link in level 1 is simply  $\pi_1$ . Thus, for links in level 1

$$\sum_{f \in \mathcal{F}_1} b(f) = \pi_1 \sigma_1 = c_1$$

So far, the design is proportional. Now assume that for all links in levels  $1 \leq i \leq N-1$ , the capacity corresponds to that of a proportional design:

$$c_i = \sum_{f \in \mathcal{F}_i} b(f)$$

Each switch in level  $i-1$  has  $n_i$  children, all with symmetric traffic. Thus, for each class of flows  $\sigma_1 \dots \sigma_{i-1}$ , the number of flows of that class that traverse each link in level  $i-1$  is  $n_i$  times the number that traverse each link in level  $i$ . So the

contribution of flows in the classes  $\sigma_1 \dots \sigma_{i-1}$  to  $\sum_{f \in \mathcal{F}_i} b(f)$  is simply  $c_{i-1}/n$ . In addition,  $\pi_i$  flows traverse each link in level  $i$  that do not traverse level  $i-1$ . These flows are of size  $\sigma_i$ . Thus for all  $l$  in level  $i$

$$\sum_{f \in \mathcal{F}_i} b(f) = \pi_i \sigma_i + \frac{c_{i-1}}{n_i} = c_i$$

By induction, the formula holds for links in all levels  $1 \leq i \leq L$ . Since the traffic pattern is interference free and the design is proportional, by Theorem 2.10 the design is optimal with respect to completion time and latency.  $\square$

### A.7 Proof of Lemma 3.3

**Design space of fat-trees**  $FT(n, 2)$ . Consider a fat-tree  $FT(n, 2)$  and let  $c_1$  and  $c_2$  be the capacity of its spine and leaf links, respectively. Without loss of generality, let  $b(f) = 1$  for interpod flows, and  $b(f) = \sigma$  for intrapod flows. Then,

- (1) If  $\sigma \geq 1$ , the optimal design satisfies  $c_1 \leq n^2 \cdot c_2 / (n+1)$ .
- (2) If  $0 < \sigma < 1$ , the proportional design satisfies  $n^2 \cdot c_2 / (n+1) < c_1 < n \cdot c_2$ , but every design wastes bandwidth.
- (3) If  $\sigma = 0$ , the optimal design corresponds to  $c_1 = n \cdot c_2$ .

**PROOF.** Applying Lemma 3.2, and substituting  $n_1 = n_2 = n$  and  $L = 2$ , the proportional design is

$$c_1 = (n-1)n^2\sigma_1$$

$$c_2 = (n-1)\sigma_2 + \frac{(n-1)n^2}{n}\sigma_1$$

Substituting  $\sigma_1 = 1$ ,  $\sigma_2 = \sigma$ ,

$$c_1 = (n-1)n^2$$

$$c_2 = (n-1)\sigma + \frac{(n-1)n^2}{n} = (n-1)(\sigma+n)$$

It follows immediately that

$$c_1 = \frac{n^2}{n+\sigma}c_2$$

The inequalities in the three cases all follow from plugging the hypotheses about  $\sigma$  into the above equation. The optimality of the proportional design in cases (1) and (3) follows from the fact that the traffic patterns are interference free in these cases. It remains only to show that when  $0 < \sigma < 1$ , every design wastes bandwidth. By Theorem 2.8, it suffices to show that the proportional design wastes bandwidth. For the proportional design, the initial fairshare allocation of each link is the quotient of its capacities and the number of flows that traverse it:

$$\frac{c_1}{(n-1)n^2} = 1$$

$$\frac{c_2}{n-1+n(n-1)} = \frac{\sigma+n}{1+n} < 1$$

Since the leaf links have the smaller initial allocation, they are the initial bottlenecks for both interpod and intrapod flows, and bandwidth is wasted at the spinelinks.  $\square$

### A.8 Proof of Lemma 4.1

**Optimal 3-level folded-Clos with radix  $k$  and skewed traffic.** Assume a traffic pattern  $b(f) = \sigma_i$ , where  $f$  is a flow that traverses a link in level  $i$ , but no link in level  $i-1$ . Assume that  $\sigma_1 = 1$  and  $\sigma_2 = \sigma_3 = \sigma$ . Then, the traffic pattern is interference-free and the oversubscription parameter  $\omega(k, \sigma) = \frac{k}{2\beta(k, \sigma)}$  corresponds to the optimal design, where

$$\beta(k, \sigma) = \left\lceil \frac{(k^4 - k^3)}{2(k^3 - k^2) + \sigma(2k^2 - 8)} \right\rceil \quad (19)$$

is the number of deployed spine blocks and  $\lceil \cdot \rceil$  is the ceiling operator. Equivalently, a Clos( $k, 3$ ) with  $\beta$  spine blocks is optimal if  $\sigma_1 = 1$  and  $\sigma_2 = \sigma_3 = \sigma(k, \beta)$ , where:

$$\sigma(k, \beta) = \begin{cases} \frac{(k^3 - k^2)(k - 2\beta)}{2\beta(k^2 - 4)}, & \text{for } 1 \leq \beta < k/2 \\ 1, & \text{for } \beta = k/2 \end{cases} \quad (20)$$

**PROOF.** From Theorem 2.10, we know that the optimal design is one that makes all flows complete at the same time in the following network configuration:

- (1) Set up a flow between every pair of hosts in the folded-Clos network;
- (2) Have each host start transmitting data to every other host according to the traffic pattern  $b(f)$ .

Consider an oversubscribed configuration of the folded-Clos. From Fig. A4-a, we know that interpod flows are bottlenecked at the spine links while intrapod flows are bottlenecked at the leaf (or edge) links. Using QTBS [27], since interpod flows are bottlenecked at the top level in the bottleneck structure, their transmission rate  $r_1$  can be obtained by dividing the capacity of a spine link with the total number of flows traversing it. It's easy to see that the total number of flows traversing a spine link when  $\beta = k/2$  (i.e., an oversubscription of 1:1) is  $\rho_1 = k^3/2 - k^2/2$ . Such number increases by a factor of  $k/(2\beta)$  when only  $\beta$  spine blocks are deployed, i.e.,  $\rho_1 = (k^3/2 - k^2/2) \cdot k/(2\beta)$ . Taking a normalized spine link capacity of  $c_1 = 1$ , this yields:

$$r_1 = \frac{c_1}{\rho_1} = \frac{1}{\rho_1} = \frac{2\beta/k}{(k^3 - k^2)/2} = \frac{4\beta}{k^4 - k^3}$$

From Fig. A4, the intrapod flows are bottlenecked at the leaf links. It's easy to see that the number of flows traversing a leaf link (the third level in the folded-Clos) is  $\rho_3 = k^3/2 - 2$ . These flows will get a fair share of the leaf link capacity minus the bandwidth taken by the interpod flows (since these also traverse the leaf links), leading to the following transmission rate for both long and short intrapod flows (taking again the normalized leaf link capacity of  $c_3 = 1$ ):

$$r_2 = r_3 = \frac{1 - r_1 \cdot \rho_1}{\rho_3 - \rho_1} = \frac{2(k - 2\beta)}{k(k^2 - 4)}$$

Using Theorem 2.10, the optimal design equalizes flow completion time, thus we have:

$$\frac{1}{r_1} = \frac{\sigma}{r_3}$$

Doing some algebraic manipulation we obtain:

$$\beta = \frac{(k^4 - k^3)}{2(k^3 - k^2) + \sigma(2k^2 - 8)}$$

Since folded-Clos is a discrete network, we need to take the smallest integer that is higher than the above expression, leading to Equation (19). Equation (20) can also be easily derived from the above expression and we leave it as an exercise to the reader.  $\square$

## A.9 Proof of Lemma 4.3

**Network completion time of a 3-level folded-Clos with radix  $k$  and skewed traffic.** Assume that every host sends  $\sigma$  bits of information to every other host located in the same pod and 1 bit of information to every other host located in a remote pod. The network completion time of a 3-level folded-Clos with radix  $k$  is:

$$\mu(\beta, \sigma) = \begin{cases} \frac{k^2 + \sigma(k^3 - k^2) - 4}{\sigma(k^4 - k^3)}, & \text{if } \sigma \leq \sigma(k, \beta) \\ \frac{2c}{4\beta c}, & \text{otherwise} \end{cases} \quad (21)$$

where  $c$  is the capacity of each switch port.

PROOF. This proof is similar to the proof of Lemma 4.1 and we omit it for the sake of brevity.  $\square$

## B FORMAL DEFINITION OF THE BOTTLENECK STRUCTURE GRAPH

At the core of QTBS lies the concept of a bottleneck. The theory builds upon a definition that captures the mathematical relationship between a link and a flow bottlenecked at it:

*Definition B.1. Bottleneck link.* Let  $\mathcal{L}$  and  $\mathcal{F}$  be the sets of links and flows in a network, respectively. Let  $c_l$  and  $r_f$  be the capacity of link  $l \in \mathcal{L}$  and the transmission rate of flow  $f \in \mathcal{F}$ , respectively. We say that flow  $f$  is bottlenecked at a link  $l \in \mathcal{L}$  if and only if flow  $f$  traverses link  $l$  and  $\partial r_f / \partial c_l \neq 0$ .

The above definition provides the starting-point connection between QTBS and the congestion control problem. Intuitively, the transmission rate  $r_f$  of a flow  $f$  depends on the capacity  $c_l$  of its bottleneck link  $l$ . This is expressed mathematically with the expression  $\partial r_f / \partial c_l \neq 0$ ; that is, a flow  $f$  is bottlenecked at a link  $l$  if and only if a change in the capacity of  $l$  affects flow  $f$ 's transmission rate. Link  $l$  here

constitutes the point of congestion of flow  $f$ . The key to a congestion control algorithm is to identify each flow's point of congestion and determine an optimal transmission rate that both maximizes throughput without creating congestion while ensuring fairness among all flows. Further, this definition allows us to introduce the concept of *bottleneck structure* [26, 27], the core building block of QTBS used to model the system-wide performance of a network:

*Definition B.2. Bottleneck Structure.* Let  $\mathcal{L}$  and  $\mathcal{F}$  be the set of links and flows in a network, respectively. The *bottleneck structure* is the directed graph such that:

- (1) There exists a vertex for each link and each flow.
- (2) If  $f \in \mathcal{F}$  traverses link  $l \in \mathcal{L}$ , then there exists a directed edge from  $f$  to  $l$ .
- (3) If  $f \in \mathcal{F}$  is bottlenecked at link  $l \in \mathcal{L}$ , then there exists a directed edge from  $l$  to  $f$ .

As shown in [26, 27], the bottleneck structure of a network describes how perturbations (small variations) in link capacities and flow transmission rates propagate through the network. Intuitively, imagine that flow  $f$  is bottlenecked at link  $l$ . From Definition B.1, this necessarily implies that a perturbation in the capacity of link  $l$  will cause a change on the transmission rate of flow  $f$ ,  $\partial r_f / \partial c_l \neq 0$ . This is reflected in the bottleneck structure by the presence of a directed edge from link  $l$  to flow  $f$  (Condition 3 in Definition B.2). A change in the value of  $r_f$ , in turn, affects all the other links traversed by flow  $f$ . This is reflected by the directed edges from  $f$  to the links it traverses (Condition 2). This process of (1) inducing a perturbation in a vertex (either in a link or a flow vertex) followed by (2) propagating the effects of the perturbation along the departing edges of the vertex creates a ripple effect in the bottleneck structure much like an instantaneous picture of a wave traveling through water, with some flows and links seeing an increment in the available bandwidth (the crest of the wave) and others seeing a reduction (the trough of the wave).

## C BOTTLENECK STRUCTURE OF FAT-TREES

For the sake of illustration, consider the case of the binary, 2-level fat-tree  $FT(2, 2)$  shown in Fig. 1. To study the behavior of this network, we build its possible bottleneck structures by using the *GradientGraph* algorithm introduced in [26, 27] and plot them in Fig. A1. Depending on the design, the interconnect has one of three possible bottleneck structures, which we call *oversubscribed* (Fig. A1-a), *balanced* (Fig. A1-b)



and *undersubscribed* (Fig. A1-c)<sup>10</sup>. Following a similar convention used in [27], each white vertex in the bottleneck structure represents a link, while each colored vertex represents a flow, with each flow’s color indicating its path according to the coloring scheme of Fig. 1. (Note that since there is bidirectional communication between every pair of hosts, for each different color there are two vertices.)

These bottleneck structures allow us to make initial qualitative observations about the design problem for fat-trees. Consider a design in which the spine links are oversubscribed; for example, imagine that the leaf links each had a capacity of 100 and the spine links each had a capacity of 1. Clearly the intrapod flows, which use the spine links, will experience a lower transmission rate than the interpod flows, which do not. This observation is reflected in the bottleneck structure corresponding to the oversubscribed configuration (Fig. A1-a). Links 1 and 2 (the spine links) are at the top of the graph, since they constrain the flows that traverse them to have a very small rate. The interpod flows are bottlenecked at the spine links, so they are directly adjacent to them in the graph. The leaf links (links 3-6) do not bottleneck the interpod flows, but they are traversed by them; so according to the definition of bottleneck structure [26], they are directly adjacent to these flows in the graph. Finally, the intrapod flows are at the bottom of the graph because they only traverse the leaf links and so they get a very fast rate. This design is desirable when the intrapod flows transmit much more traffic than the interpod flows, because the intrapod flows get a correspondingly faster rate. But if all flows transmit the same amount of data, then the intrapod flows will finish transmitting long before the interpod flows, leading to wasted bandwidth at the leaf links later on.

Next consider a design in which the spine links are undersubscribed; for example, one where the leaf links have capacity of 1 and the spine links have capacity of 100. Now all flows (interpod and intrapod) are bottlenecked by the leaf links and experience equal rates, and the spine links aren’t bottlenecks at all. This is reflected in the bottleneck structure for this case (Fig. A1-c), where all flow vertices lie below the leaf links, and the spine links have no children in the graph. This bottleneck structure shows that this design wastes bandwidth, since we can reduce the capacity (and the cost) of the spine links without harming the performance of any flows. The balanced configuration is less costly (as it requires less capacity in the spine links) and is as performant as the undersubscribed configuration (since in both configurations all flows are bottlenecked at the leaf links and, thus, will experience the same throughput). Another way of

seeing this is to note that spine links are not bottlenecks for any flows and, thus, some of their capacity is being wasted.

Finally, for the balanced configuration (Fig. A1-b), all flows are bottlenecked at both the spine and the leaf links equally. Since all flows are on the same level of the bottleneck structure, they will all get the same rate. In Section 3.1 we prove that this design is optimal when the traffic pattern is uniform (Definition 2.2). Since production networks typically experience skewed traffic patterns, with intrapod flows sending more data than interpod flows [4, 5, 17], this design is usually sub-optimal.

In the next section we introduce the set of optimal designs and mathematically characterize some of these properties inherent to fat-trees.

## D TRAFFIC SKEWNESS AND INTERCONNECT SIZE FOR $FT(n, 2)$

In Fig. A2 we present a chart representing the traffic skewness value needed to optimally operate a  $FT(n, 2)$  as a function of the number of hosts  $n^2$  supported by the interconnect and for various tapering parameters  $\tau$ . As shown, the higher the degree of oversubscription (lower value of  $\tau$ ), the higher the skewness level required to operate efficiently. In the limit where there is no oversubscription ( $\tau = 1$ ), the optimal skewness is a horizontal line  $\sigma = 1$  corresponding to the case of uniform traffic and a balanced bottleneck structure (e.g., Fig. A1-b).

## E NETWORK COMPLETION TIME OF $Clos(48, 3)$

Fig. A3 provides a plot of the network completion time as a function of the oversubscription parameter  $\omega$  for a production-scale folded-Clos with radix  $k = 48$  assuming a normalized link capacity of  $c = 1$  bps. The plot is based on the network completion time Equation (13). In this chart, network completion time has been normalized to (divided by)  $\sigma$  to better illustrate its asymptotic behavior, so the corresponding traffic pattern is  $b(f) = 1/\sigma$  for interpod flows and  $b(f) = 1$  for intrapod flows.

## F BOTTLENECK STRUCTURE OF FOLDED-CLOS NETWORKS

Fig. A4 provides the possible bottleneck structures of the  $Clos(4, 3)$  interconnect. Link labels inside the white vertices and flow colors correspond to those used in Fig. 4. Because folded-Clos are discrete networks, there are only two possible bottleneck structures: oversubscribed (Fig. A4-a) and undersubscribed (Fig. A4-b). (Thus, unlike fat-trees, folded-Clos don’t have a balanced bottleneck structure).

<sup>10</sup>Note that when we use the terms oversubscribed and undersubscribed, we mean with respect to the balanced solution—not with respect to a full fat-tree, like some previous authors [1].

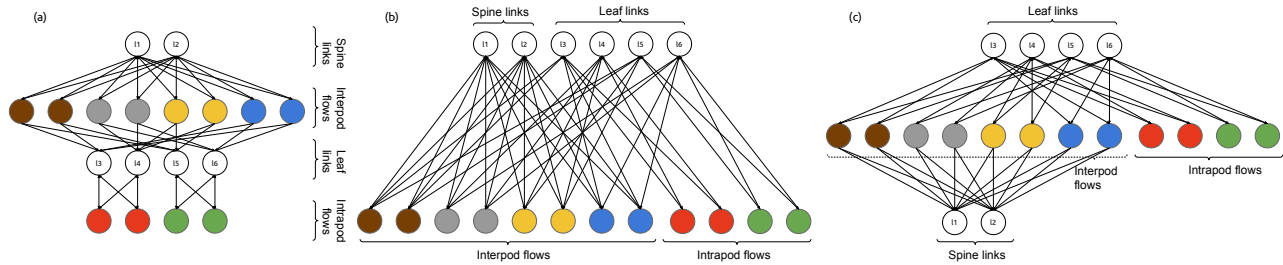


Figure A1: Possible bottleneck structures in a  $FT(2, 2)$ : (a) oversubscribed, (b) balanced, (c) undersubscribed

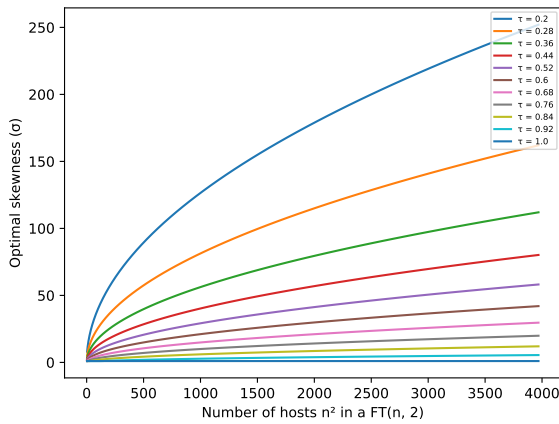


Figure A2: Optimal traffic skewness in  $FT(n, 2)$  as a function of the number of hosts  $n^2$ .

## G DESIGNING DRAGONFLY NETWORKS

While folded-Clos is the predominant topology in large-scale data centers, evolving technology and the availability of high-radix switches have led to other new high-performance topologies that, for some applications, are able to deliver better cost-performance trade-off. One such topology is dragonfly, introduced by Kim et al. [18], which leverages modern switches with high radix to reduce the diameter, the latency and in some cases the cost of the interconnect. Dragonfly topologies are primarily being used in supercomputer interconnects (e.g., [29]).

A dragonfly interconnect consists of  $p$  pods<sup>11</sup>, each with  $a$  switches. Each switch is connected with every other switch in its pod via intrapod links, forming a full-mesh. Each switch is also connected with  $h$  other switches located in other pods via interpod links. Finally, each switch is connected with  $t$  hosts. The switches in a dragonfly must offer at least a

<sup>11</sup>In the context of supercomputer interconnects, these are usually referred as *groups*. We use the data center-oriented term *pod* for consistency with the terminology used in the analysis of the previous interconnects.

radix  $k = t + h + a - 1$  and the interconnect scales to support  $a \cdot p \cdot t$  hosts. We will use the notation  $Dragonfly(a, p, h)$  to denote a dragonfly with parameters  $a, p$  and  $h$ . (For the sake of simplicity, we omit the connections with the hosts, since they do not alter the core topology of the network.) A dragonfly is said to be *canonical* if  $p = a + 1$  and  $h = 1$ , denoted as  $Dragonfly(a, a + 1, 1)$ . Canonical dragonflies are of interest because they ensure every pair of hosts can be connected by traversing one single interpod link and have minimal diameter.

The topology of a dragonfly interconnect is shown in Fig. A5 for the case of  $a = 3, p = 14$  and  $h = 1$ . A dragonfly can generally be constructed using the following iterative procedure. Position all the pods in a circle as shown in the figure and start with an arbitrary pod. Refer to this single pod as *group 1*. From it, connect its  $a \cdot h$  interpod links to as many consecutive different pods as possible by traveling the ring counterclockwise. Refer to this set of pods *group 2*. Now, connect its  $ah(ah - 1)$  interpod links to as many consecutive different pods as possible by continuing to travel the ring counterclockwise. Refer to this set of pods *group 3*. This process is repeated until all the interpod links have been connected. At the end of this process, we set  $\lambda$  to be the total number of groups and  $\gamma$  to be the number of pods in the last group. In Section G.2 we will see that the parameters  $\lambda$  and  $\gamma$ , together with  $a, p$  and  $h$ , uniquely determine the performance of a dragonfly interconnect.

### G.1 Bottleneck Structure of Dragonflies

Fig. A6 provides the possible bottleneck structures of the  $Dragonfly(3, 4, 1)$  interconnect. Flow colors correspond to those used in Fig. A5. There are three possible bottleneck structures: oversubscribed (Fig. A6-a), balanced (Fig. A6-b) and undersubscribed (Fig. Fig. A6-c).

### G.2 Design Equations for Skewed Traffic

The general equations that determine an optimal dragonfly design for skewed traffic are as follows:

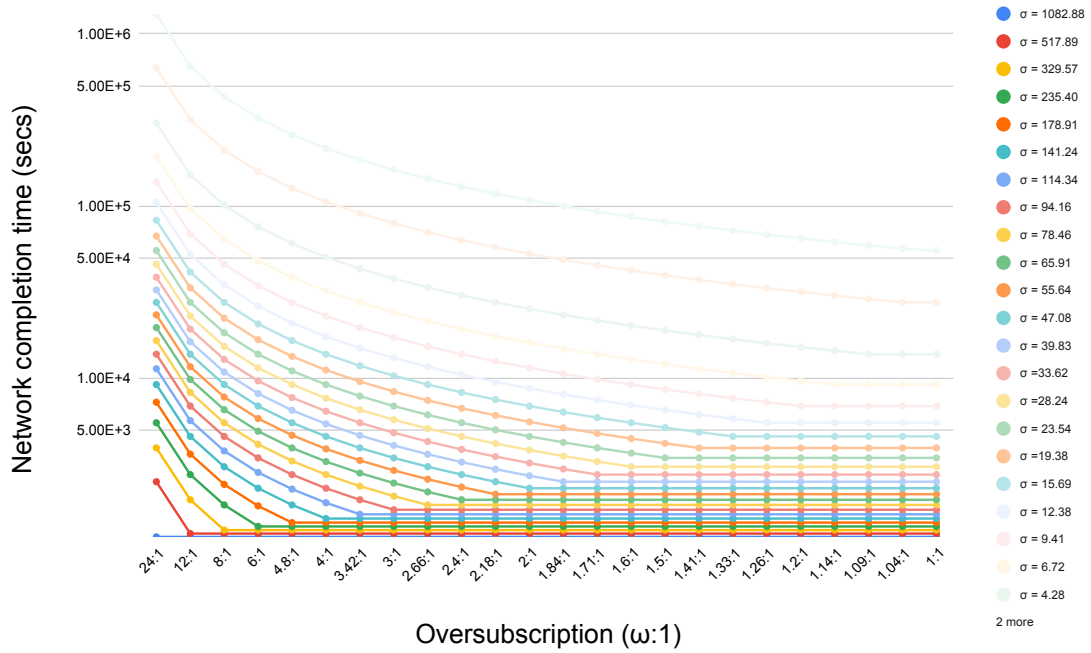


Figure A3: Effect of oversubscribing a Clos(48, 3) network on the maximum flow completion time.

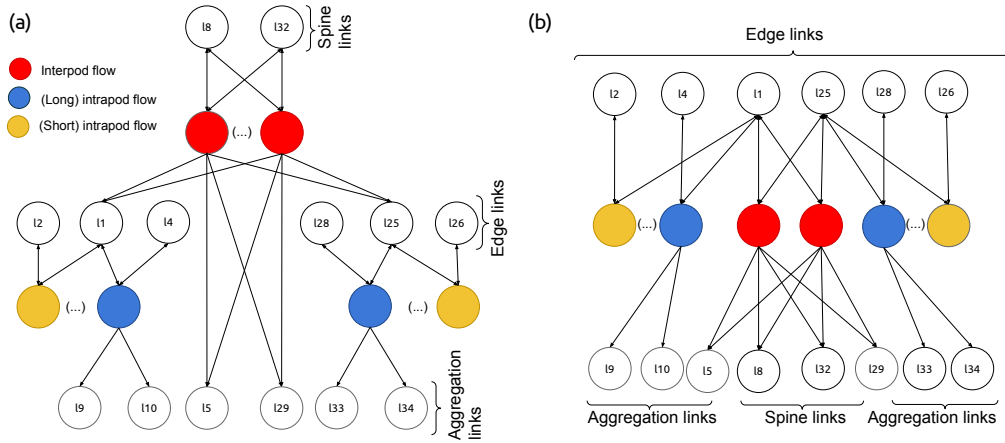


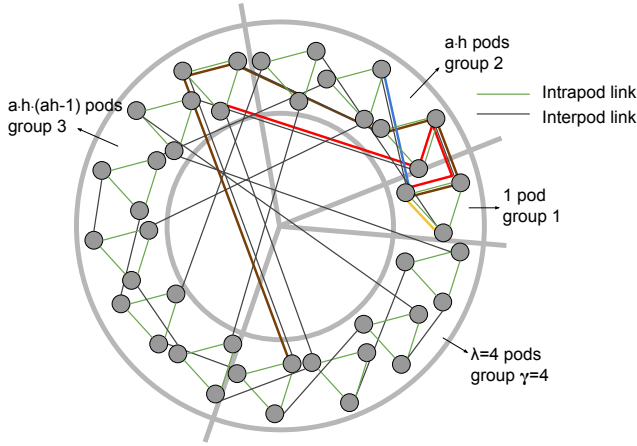
Figure A4: Possible bottleneck structures in a Clos(4, 3): (a) oversubscribed and (b) undersubscribed

LEMMA G.1. *Optimal dragonfly with skewed traffic. Assume that every host sends  $\sigma$  bits of information to every other host located in the same pod and 1 bit of information to every other host located in a remote pod. If  $\sigma \geq 1$ , then the traffic pattern is interference-free and the following design is optimal:*

$$c_1 = \frac{\rho_1}{\sigma + \rho_2 - 1} c_2 \tag{22}$$

where

$$\rho_1 = \frac{1}{h} \left[ a(g-1) + \sum_{i=0}^{y-2} \left( a(g-1) - \sum_{j=0}^i a^2 h (ah-1)^j \right) \right] \tag{23}$$



**Figure A5: Dragonfly topology with  $a = 3, p = 14, h = 1$ , resulting in  $\lambda = 4$  and  $\gamma = 4$ .**

$$\rho_2 = g + \lambda + \frac{1}{a-1} \sum_{i=0}^{\gamma-2} \left[ ah(a-1)(ah-1)^i + \frac{(a-1)h}{ah-1} \left( a(g-1) - \sum_{j=0}^i a^2 h (ah-1)^j \right) \right] \quad (24)$$

$$\gamma = \max \left\{ x \mid a(g-1) - \sum_{i=0}^{x-2} a^2 h (ah-1)^i \geq 0 \right\} \quad (25)$$

$$\lambda = g - 1 - \sum_{i=0}^{\gamma-1} ha(ah-1)^i \quad (26)$$

PROOF. This proof is similar to the proof of Lemma 3.2 and we omit it for the sake of brevity.  $\square$

With simple algebraic manipulations, we can use the above lemma to derive the optimal design for a canonical dragonfly by setting  $p = a + 1$  and  $h = 1$ :

$$c_1 = \frac{a^2}{\sigma + 2a} c_2 \quad (27)$$

Since a canonical dragonfly has a total of  $a(a+1)/2$  interpod links and  $a(a-1)(a+1)/2$  intrapod links, its tapering parameter  $\tau$  corresponds to<sup>12</sup>:

$$\tau = \frac{c_1 \cdot a \cdot (a+1)/2}{c_2 \cdot a \cdot (a-1)(a+1)/2} = \frac{a^2}{(\sigma + 2a)(a-1)} \quad (28)$$

For  $\sigma = 0$ , we have that the optimal design has a tapering parameter  $\tau = a/(2(a-1))$ . Interestingly, this is in contrast with the optimal fat-tree design for  $\sigma = 0$ , which corresponds to a full fat-tree with  $\tau = 1$ . This implies that dragonflies require a maximum tapering parameter which is lower than

<sup>12</sup>Recall from Section 3 that we defined  $\tau$  as the ratio of the aggregated interpod link capacity divided by the aggregated intrapod link capacity

that of a fat-tree. In other words, an optimal design for  $\sigma = 0$  requires less aggregated capacity in the interpod links relative to the intrapod links for canonical dragonflies than it requires for fat-trees. This also implies that designing a dragonfly by allocating as much aggregated capacity to the interpod links as the intrapod links (i.e.,  $\tau = 1$  and  $c_1 = (a-1) \cdot c_2$ ) is always suboptimal, as the optimal design for  $\sigma = 0$  (i.e.,  $c_1 = a \cdot c_2/2$ ) is less costly and yields the same network completion time and throughput.

Fig. A7 shows the optimal tapering parameter  $\tau$  as a function of traffic skewness  $\sigma$  for a variety of canonical dragonflies. Similar to the case of fat-trees in Fig. 3, all optimal designs require oversubscribing the interpod links ( $\tau < 1$ ) and, as the size of the interconnect increases, the optimal tapering parameter increases too. As an example, suppose that our goal is to design a canonical dragonfly  $Dragonfly(6, 7, 1)$  (i.e.,  $a = 6$ ) to transport a traffic pattern with skewness  $\sigma = 15$ . By using the chart in Fig. 3, we can identify the needed design (represented with a yellow dot) at the intersection of the red curve with the line  $\sigma = 15$ , resulting in a tapering parameter of  $\tau = .26667$ .

The following lemma characterizes the design space of canonical dragonflies  $Dragonfly(a, a+1, 1)$ :<sup>13</sup>

LEMMA G.2. *Design space of canonical dragonflies. Consider a canonical dragonfly  $Dragonfly(a, a+1, 1)$  and let  $c_1$  and  $c_2$  be the capacity of its interpod and intrapod links, respectively. Without loss of generality, let  $b(f) = 1$  for interpod flows, and  $b(f) = \sigma$  for intrapod flows. Then,*

- (1) If  $\sigma \geq 1$ , the optimal design satisfies  $c_1 \leq a^2 \cdot c_2 / (1 + 2a)$ .
- (2) If  $\sigma < 1$ , the optimal design satisfies  $a^2 \cdot c_2 / (1 + 2a) < c_1 \leq a \cdot c_2 / 2$ .

PROOF. See appendix A.7.  $\square$

The above lemma is pictured in Fig. A8 as follows. The region of optimal designs for  $\sigma > 1$  is marked with gray hash lines; it falls below the optimal design line for uniform ( $\sigma = 1$ ) traffic  $c_1 = a^2 \cdot c_2 / (1 + 2a)$ , shown as a red line. Note that the optimal designs for  $\sigma > 1$  correspond to oversubscribed interconnects (shown in Fig. A6-a) because they allocate less bandwidth to the spine links than a balanced design. The region of optimal designs for  $\sigma < 1$  is marked with red hash lines, bordered by the design  $c_1 = a \cdot c_2 / 2$  (corresponding to the optimal design when  $\sigma = 0$ ) and the optimal design line for uniform traffic. These designs correspond to under-subscribed interconnects (shown in Fig. A1-c) because they allocate more bandwidth to the spine links than a balanced design. Similarly to fat-trees (although bordered by a different design curve), there exists no traffic pattern for which a design in the region  $c_1 > a \cdot c_2 / 2$  is optimal, because the

<sup>13</sup>This lemma can be generalized to support  $Dragonfly(a, p, h)$ , we leave this as an exercise to the reader.

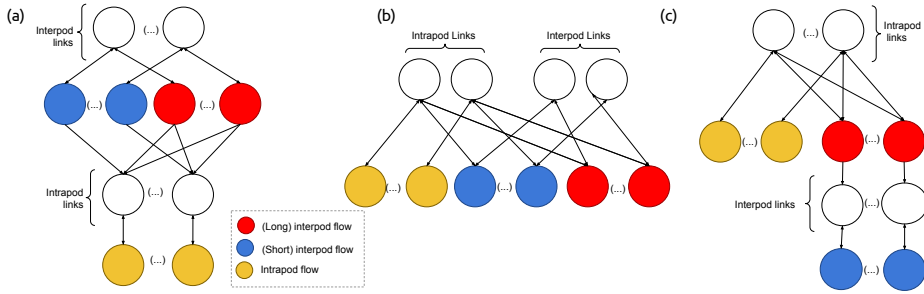


Figure A6: Possible bottleneck structures in a *Dragonfly*(3, 4, 1): (a) oversubscribed, (b) balanced and (c) undersubscribed

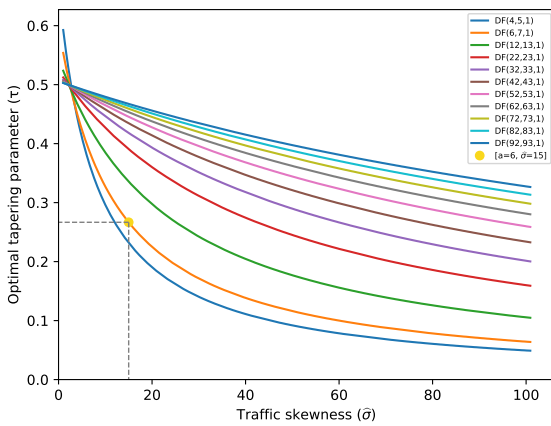


Figure A7: Optimal tapering parameter as a function of traffic skewness for various *Dragonfly*( $a, a + 1, 1$ ) designs.

design  $c_1 = a \cdot c_2/2$  yields the same network completion time and throughput and is less costly. Such designs are *inefficient* and network architects should avoid using them no matter the traffic pattern.

As an example, a *Dragonfly*(16, 17, 1) design in which intrapod flows carry ten times more traffic than interpod flows ( $\sigma = 10$ ) has an optimal tapering parameter  $\tau = 0.40635$ , yielding the design  $c_1/c_2 = 6.09523$ . This design is shown in Fig. A8 as a blue dot. Similarly, the same interconnect with uniform traffic ( $\sigma = 1$ ) has an optimal tapering parameter  $\tau = 0.51717$ , yielding the design  $c_1/c_2 = 7.75757$ . This corresponds to an interconnect with a balanced bottleneck structure (e.g., Fig. A6-b), and is plotted in Fig. A8 as a green dot.

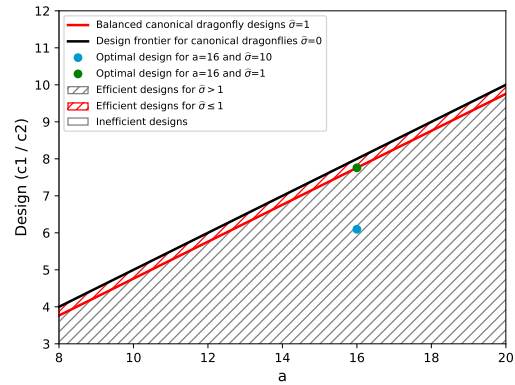


Figure A8: Design space for *Dragonfly*( $a, a + 1, 1$ ).

### G.3 Experiments with Dragonfly Networks

In this section we empirically demonstrate the existence of optimal dragonfly designs. (We follow a similar methodology used in Section 5.1 to empirically validate optimal fat-tree designs.) We start by simulating a *Dragonfly*(3, 4, 1) interconnect—i.e., a dragonfly with 4 pods, each with 3 routers (a total of 12 routers), every router in a pod is connected with every other router in the same pod (full-mesh connectivity) using intrapod links (a total of 12 intrapod links) and each router is connected with a router located in a different pod via an interpod link (a total of 6 interpod links). We connect every router to a host (a total of 12 hosts) and connect every pair of nodes with two TCP flows (one for each direction), for a total of 132 flows. In the first set of experiments, we assume uniform traffic ( $\sigma = 1$ ). Using Equation (27), we have that the optimal design corresponds to  $c_1 = a^2 \cdot c_2 / (1 + 2a) = 3^2 \cdot c_2 / (1 + 2 \cdot 3) = 1.2857 \cdot c_2$ .

Fig. A9 shows the result of simulating a variety of designs with  $c_2 = 20$  Mbps and  $c_1 \in \{5, 10, 20, 25.71, 30, 40, 60\}$  Mbps—resulting in values for  $c_1/c_2$  of 0.25, 0.5, 1, 1.2857, 1.5, 2 and 3. (Again, without loss of generality, we could pick any value for



$c_2$  and scale  $c_1$  accordingly). Results are shown for the BBR, Cubic and Reno congestion control algorithms, and for both experimental and theoretical (according to QTBS) values. As predicted by the QTBS mathematical model, the plots show that the optimal design is found at  $c_1/c_2 = 1.2857$ . According to Theorem 2.10, this design wastes no bandwidth, minimizes network completion time and maximizes network throughput. Designs in the region  $c_1/c_2 < 1.2857$  waste bandwidth at the intrapod links, increase network completion time (due to the longer completion time of the interpod flows) and decrease network throughput. As shown also, a design in the region  $c_1/c_2 > 1.2857$  achieves the same network completion time and network throughput as the design  $c_1/c_2 = 1.2857$ , regardless of how large the capacity of a spine link ( $c_1$ ) is. A design in this region wastes bandwidth at the spine links and is more costly than the optimal design, thus it should also be avoided.

The three congestion control algorithms closely follow the theoretical QTBS model, all reaching an inflection point right at the optimal design  $c_1/c_2 = 1.2857$ . Increasing the capacity of  $c_1/c_2$  beyond this value yields a design with the same network's completion time (and, thus the same network throughput) but with a higher cost. Once again and similar to the results shown in Section 5.1, BBR is able to more accurately perform according to the optimal design than Cubic and Reno. This is also shown in Fig. A10, where BBR yields a higher Jain's index. Note that, as in all other experiments, all the results follow the model for network completion time from slightly above because of imperfections of the congestion control algorithms.

Fig. A11 shows the cumulative distribution function of all the experiments run in Fig. A9a. The figure shows how increasing  $c_1$  helps reduce the maximum flow completion time until the optimal value  $c_1 = 1.2857 \cdot c_2$  is reached. At this point, all completion times are equalized (Theorem 2.10) and increasing  $c_1$  beyond this value does not qualitatively alter the completion time of the flows.

In Fig. A12 we present experiments using a skewed traffic pattern consisting of  $b(f) = 1$  for interpod flows and  $b(f) = \sigma$  for intrapod flows, with  $\sigma \in \{2, 4, 10\}$ . Using Equation (27), we have that the optimal design satisfies  $c_1 = a^2 \cdot c_2 / (\sigma + 2a) = 9 \cdot c_2 / (\sigma + 6)$ . This leads to three optimal designs, one for each traffic pattern:  $c_1 = 1.125 \cdot c_2$  for  $\sigma = 2$ ,  $c_1 = 0.9 \cdot c_2$  for  $\sigma = 4$ , and  $c_1 = 0.5625 \cdot c_2$  for  $\sigma = 10$ . Fig. A12 shows that, experimentally, the interconnect behaves as QTBS predicts, minimizing network completion time right at the optimal design for each of traffic skewness value. As predicted by the model, for any design allocating more capacity in the interpod links than the optimal design, network completion time does not improve. Network architects should avoid these designs.

## H G2-MININET

The G2-Mininet tool [19] provides a powerful, flexible interface to emulate networks of choice with customizable topology, routing and traffic flow configurations, with a focus to help experimentally demonstrate the quantitative theory of bottleneck structures (QTBS). It uses Mininet [22] and the POX SDN controller [25] to create such highly customizable networks. It also uses iPerf [13] internally to generate network traffic and offers an interface to configure various flow parameters such as the source and destination hosts, routing, start time, data size, and traffic pattern, among others. G2-Mininet also offers an integration with sFlow-RT [24] agent that enables real-time access to traffic flows and real-time computation of the emulated network's bottleneck structure. The extensions to Mininet include scripts to automatically generate the specifications of data center networks such as fat-trees, folded-Clos and dragonflies, that are then used as inputs to the emulation environment to measure the accuracy of the QTBS model and equations. Mininet uses real, production grade TCP/IP stack from the Linux kernel, enabling a testbed to run experiments using congestion control protocols such as BBR, Cubic and Reno to study bottleneck structures and flow performance. Apart from its flexible configuration interface, G2-Mininet also offers a set of useful utilities to compute and plot various performance metrics such as instantaneous network throughput, flow convergence time, flow completion time, Jain's fairness index, and the computation in real time of the network's bottleneck structure, among others for a given experiment. *G2-Mininet* and all the experiments presented in this paper are made available as open source software (see [19]).

## I ADDITIONAL EXPERIMENTS

As mentioned in Section 5, more than 600 simulation (or the equivalent of 800 hours) were run to verify that the three interconnects behave according to the design equations presented in this paper. Because only a limited selection of these experimental results are included in the main body of this paper, in Figures A13, A14, A15, A16, A17, A18, A19 and A20 we include other complementary results.

## J DESIGN TABLES FOR $FT([n_1, n_2])$

This section presents sample tables that can help network designers identify optimal designs as a function of network size and traffic pattern. Tables can be derived for any of the studied interconnects using the design equations provided in this paper. For the sake of illustration, we only provide two examples for the design of fat-trees for skewness  $\sigma = 1$  (Table 1) and  $\sigma = 2$  (Table 2). The parameters shown in this table are:

- $c_1$ : capacity of spine links.

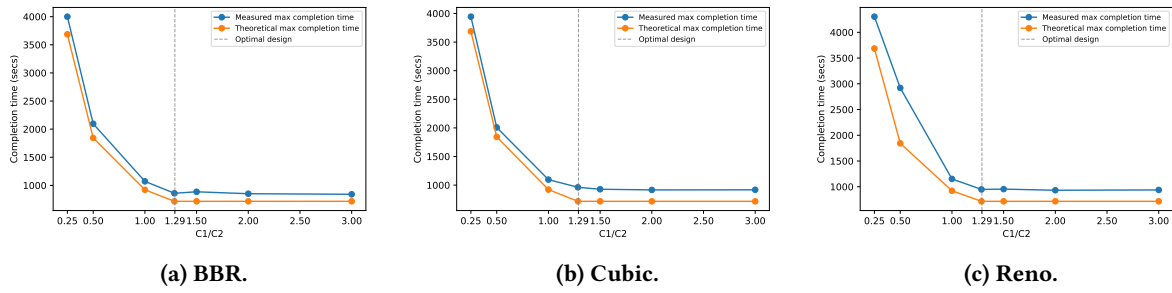
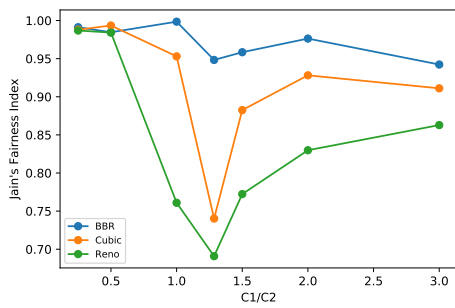


Figure A9: Network completion time for *Dragonfly*(3, 4, 1) and uniform traffic.



[H]

Figure A10: Jain's fairness index for *Dragonfly*(3, 4, 1).

- $c_2$ : capacity of leaf links.
- $r_1$ : transmission rate of interpod flows.
- $r_2$ : transmission rate of intrapod flows.
- $s_1$ : fair share of spine links.
- $s_2$ : fair share of leaf links.
- $\tau$ : tapering parameter.
- $\eta$ : total flow throughput divided by total link capacity.<sup>14</sup>

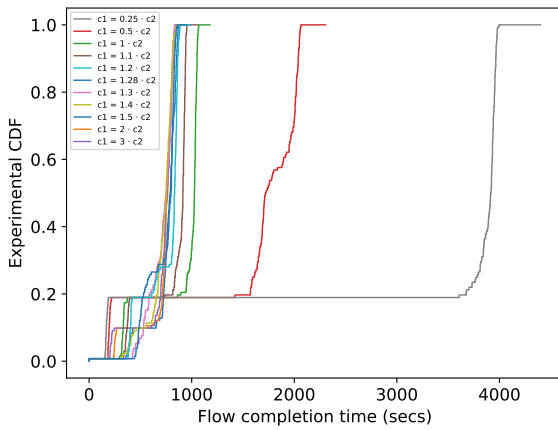


Figure A11: CDFs of flow completion time for *Dragonfly*(3, 4, 1).

<sup>14</sup>This metric provides a simple estimate of how efficient is this design.

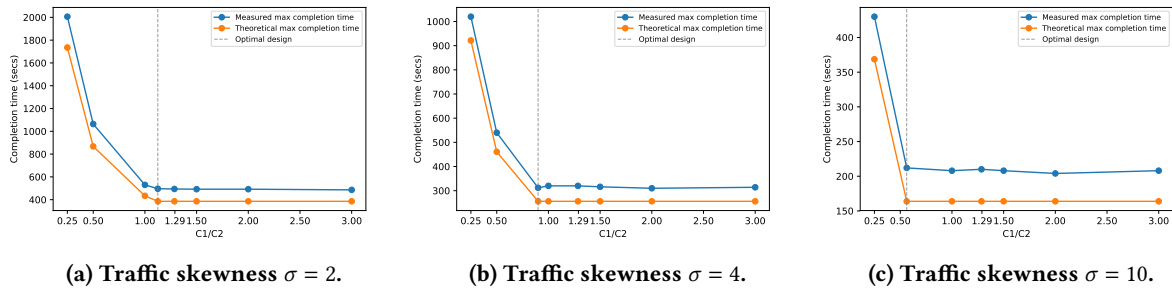


Figure A12: Network completion time for *Dragonfly*(3, 4, 1) using BBR.

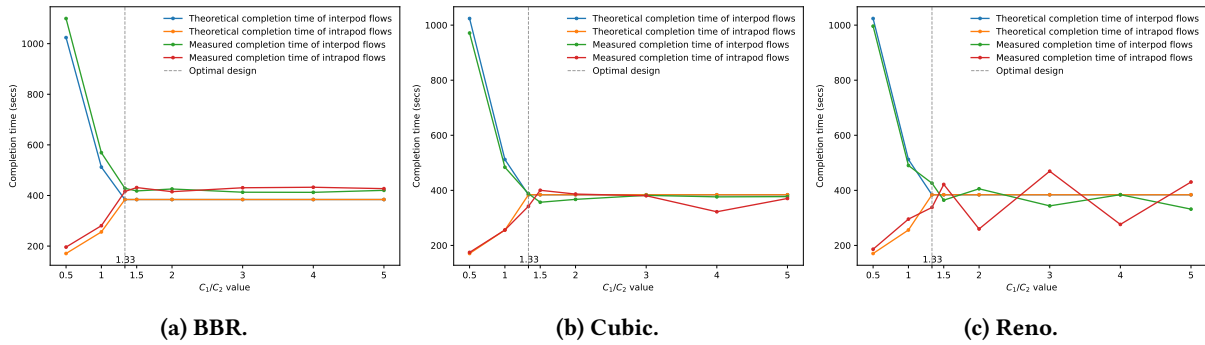


Figure A13: Flow completion time for *FT*(2, 2) and uniform traffic.

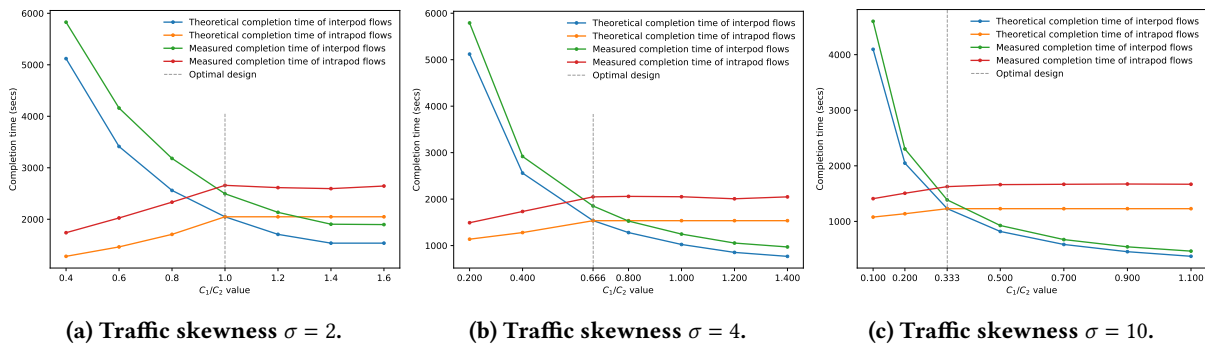


Figure A14: Flow completion time for *FT*(2, 2) using BBR.

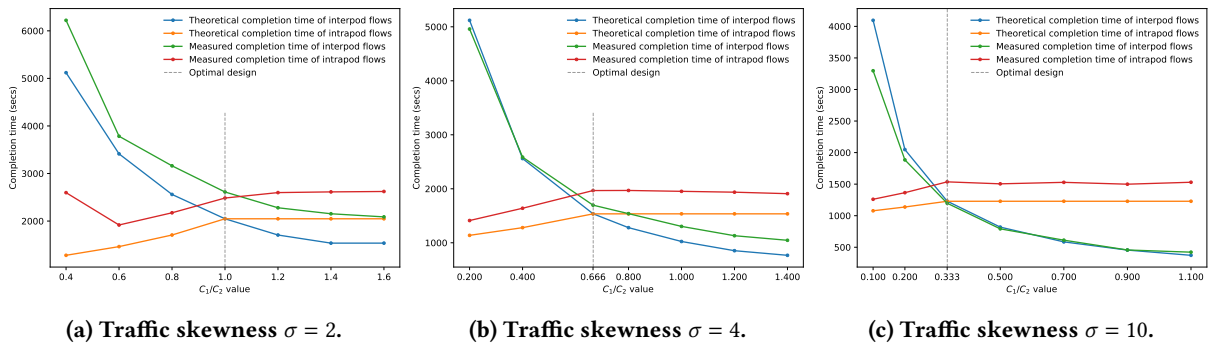


Figure A15: Flow completion time for  $FT(2, 2)$  using Cubic.

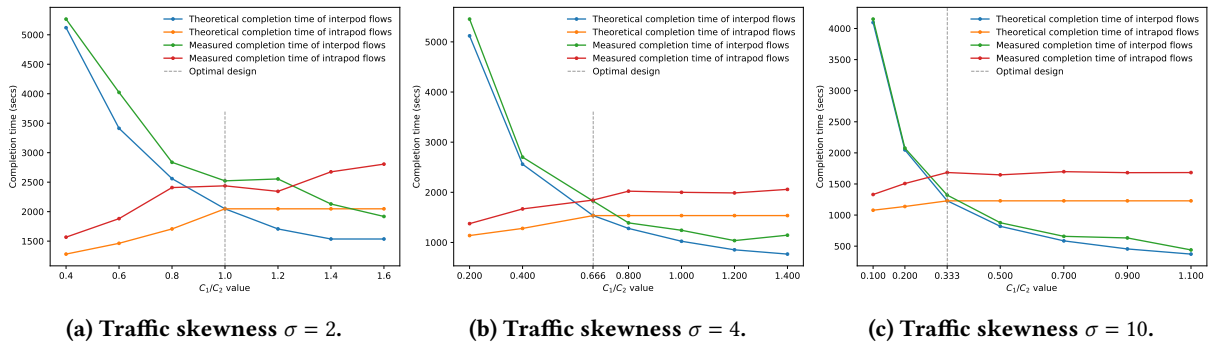


Figure A16: Flow completion time for  $FT(2, 2)$  using Reno.

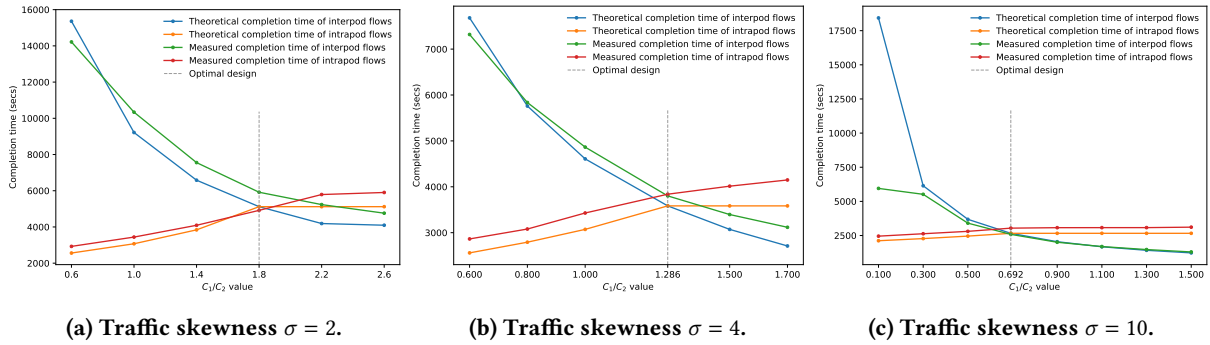


Figure A17: Flow completion time for  $FT(3, 2)$  using Cubic.

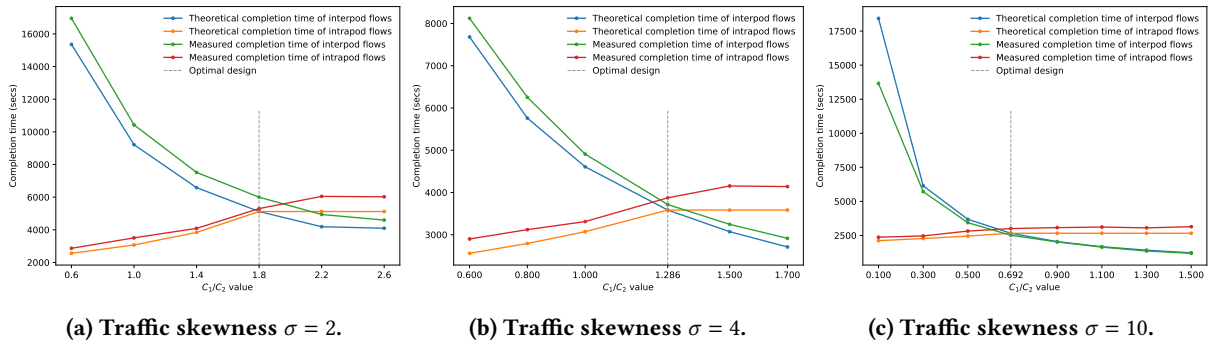


Figure A18: Flow completion time for  $FT(3, 2)$  using Reno.

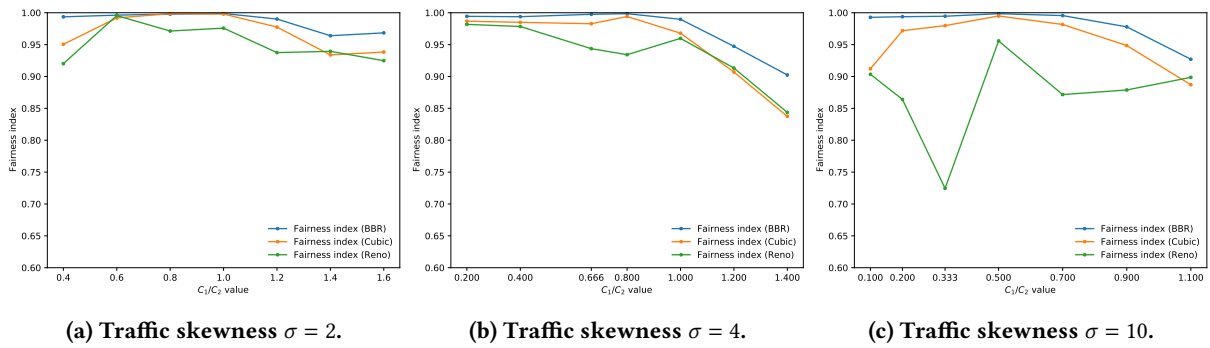


Figure A19: Jain's fairness index for  $FT(2, 2)$ .

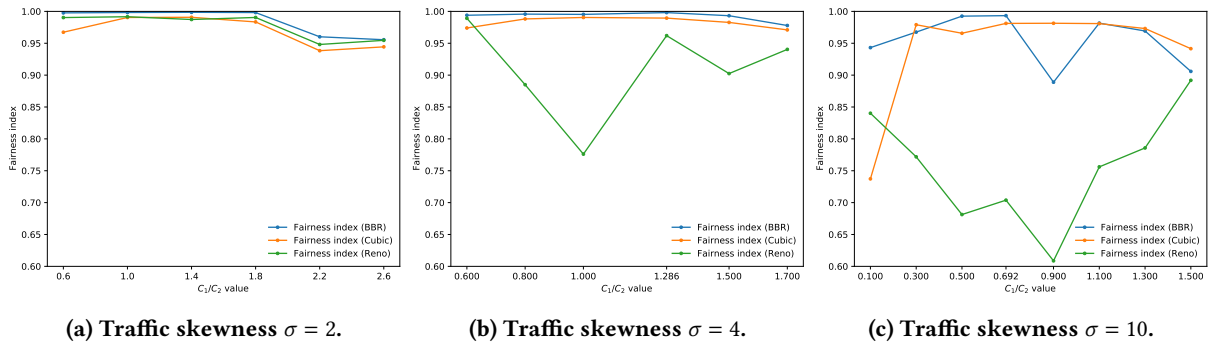


Figure A20: Jain's fairness index for  $FT(3, 2)$ .



**Table 1: Optimal fat-tree designs  $FT([n_1, n_2])$  for skewness  $\sigma = 1$  (uniform traffic).**

$n_1$	$n_2$	$c_1/c_2$	$s_1$	$s_2$	$\eta$	$\tau$	$n_1$	$n_2$	$c_1/c_2$	$s_1$	$s_2$	$\eta$	$\tau$	$n_1$	$n_2$	$c_1/c_2$	$s_1$	$s_2$	$\eta$	$\tau$
2	4	1.3333	0.25	0.25	0.6	1.5	3	42	9.561	0.0026	0.0026	0.5942	1.4643	8	24	2.7391	0.0159	0.0159	0.5227	1.0952
2	6	1.8	0.1111	0.1111	0.625	1.6667	3	45	10.2273	0.0022	0.0022	0.5946	1.4667	8	32	3.6129	0.0089	0.0089	0.5254	1.1071
2	8	2.2857	0.0625	0.0625	0.6364	1.75	3	48	10.8936	0.002	0.002	0.5949	1.4688	8	40	4.4872	0.0057	0.0057	0.527	1.1143
2	10	2.7778	0.04	0.04	0.6429	1.8	4	8	1.7143	0.0833	0.0833	0.5385	1.1667	8	48	5.3617	0.004	0.004	0.5281	1.119
2	12	3.2727	0.0278	0.0278	0.6471	1.8333	4	12	2.4545	0.037	0.037	0.55	1.2222	9	18	1.8824	0.0312	0.0312	0.5152	1.0625
2	14	3.7692	0.0204	0.0204	0.65	1.8571	4	16	3.2	0.0208	0.0208	0.5556	1.25	9	27	2.7692	0.0139	0.0139	0.52	1.0833
2	16	4.2667	0.0156	0.0156	0.6522	1.875	4	20	3.9474	0.0133	0.0133	0.5588	1.2667	9	36	3.6571	0.0078	0.0078	0.5224	1.0938
2	18	4.7647	0.0123	0.0123	0.6538	1.8889	4	24	4.6957	0.0093	0.0093	0.561	1.2778	9	45	4.5455	0.005	0.005	0.5238	1.1
2	20	5.2632	0.01	0.01	0.6552	1.9	4	28	5.4444	0.0068	0.0068	0.5625	1.2857	10	20	1.8947	0.0278	0.0278	0.5135	1.0556
2	22	5.7619	0.0083	0.0083	0.6562	1.9091	4	32	6.1935	0.0052	0.0052	0.5636	1.2917	10	30	2.7931	0.0123	0.0123	0.5179	1.0741
2	24	6.2609	0.0069	0.0069	0.6571	1.9167	4	36	6.9429	0.0041	0.0041	0.5645	1.2963	10	40	3.6923	0.0069	0.0069	0.52	1.0833
2	26	6.76	0.0059	0.0059	0.6579	1.9231	4	40	7.6923	0.0033	0.0033	0.5652	1.3	11	22	1.9048	0.025	0.025	0.5122	1.05
2	28	7.2593	0.0051	0.0051	0.6585	1.9286	4	44	8.4419	0.0028	0.0028	0.5658	1.303	11	33	2.8125	0.0111	0.0111	0.5161	1.0667
2	30	7.7586	0.0044	0.0044	0.6591	1.9333	4	48	9.1915	0.0023	0.0023	0.5663	1.3056	11	44	3.7209	0.0063	0.0062	0.5181	1.075
2	32	8.2581	0.0039	0.0039	0.6596	1.9375	5	10	1.7778	0.0625	0.0625	0.5294	1.125	12	24	1.913	0.0227	0.0227	0.5111	1.0455
2	34	8.7576	0.0035	0.0035	0.66	1.9412	5	15	2.5714	0.0278	0.0278	0.5385	1.1667	12	36	2.8286	0.0101	0.0101	0.5147	1.0606
2	36	9.2571	0.0031	0.0031	0.6604	1.9444	5	20	3.3684	0.0156	0.0156	0.5429	1.1875	12	48	3.7447	0.0057	0.0057	0.5165	1.0682
2	38	9.7568	0.0028	0.0028	0.6607	1.9474	5	25	4.1667	0.01	0.01	0.5455	1.2	13	26	1.92	0.0208	0.0208	0.5102	1.0417
2	40	10.2564	0.0025	0.0025	0.661	1.95	5	30	4.9655	0.0069	0.0069	0.5472	1.2083	13	39	2.8421	0.0093	0.0093	0.5135	1.0556
2	42	10.7561	0.0023	0.0023	0.6613	1.9524	5	35	5.7647	0.0051	0.0051	0.5484	1.2143	14	28	1.9259	0.0192	0.0192	0.5094	1.0385
2	44	11.2558	0.0021	0.0021	0.6615	1.9545	5	40	6.5641	0.0039	0.0039	0.5493	1.2188	14	42	2.8537	0.0085	0.0085	0.5125	1.0513
2	46	11.7556	0.0019	0.0019	0.6618	1.9565	5	45	7.3636	0.0031	0.0031	0.55	1.2222	15	30	1.931	0.0179	0.0179	0.5088	1.0357
2	48	12.2553	0.0017	0.0017	0.662	1.9583	6	12	1.8182	0.05	0.05	0.5238	1.1	15	45	2.8636	0.0079	0.0079	0.5116	1.0476
3	6	1.6	0.125	0.125	0.5556	1.25	6	18	2.6471	0.0222	0.0222	0.5312	1.1333	16	32	1.9355	0.0167	0.0167	0.5082	1.0333
3	9	2.25	0.0556	0.0556	0.5714	1.3333	6	24	3.4783	0.0125	0.0125	0.5349	1.15	16	48	2.8723	0.0074	0.0074	0.5109	1.0444
3	12	2.9091	0.0312	0.0312	0.5789	1.375	6	30	4.3103	0.008	0.008	0.537	1.16	17	34	1.9394	0.0156	0.0156	0.5077	1.0312
3	15	3.5714	0.02	0.02	0.5833	1.4	6	36	5.1429	0.0056	0.0056	0.5385	1.1667	18	36	1.9429	0.0147	0.0147	0.5072	1.0294
3	18	4.2353	0.0139	0.0139	0.5862	1.4167	6	42	5.9756	0.0041	0.0041	0.5395	1.1714	19	38	1.9459	0.0139	0.0139	0.5068	1.0278
3	21	4.9	0.0102	0.0102	0.5882	1.4286	6	48	6.8085	0.0031	0.0031	0.5402	1.175	20	40	1.9487	0.0132	0.0132	0.5065	1.0263
3	24	5.5652	0.0078	0.0078	0.5897	1.4375	7	14	1.8462	0.0417	0.0417	0.52	1.0833	21	42	1.9512	0.0125	0.0125	0.5062	1.025
3	27	6.2308	0.0062	0.0062	0.5909	1.4444	7	21	2.7	0.0185	0.0185	0.5263	1.1111	22	44	1.9535	0.0119	0.0119	0.5059	1.0238
3	30	6.8966	0.005	0.005	0.5918	1.45	7	28	3.5556	0.0104	0.0104	0.5294	1.125	23	46	1.9556	0.0114	0.0114	0.5056	1.0227
3	33	7.5625	0.0041	0.0041	0.5926	1.4545	7	35	4.4118	0.0067	0.0067	0.5312	1.1333	24	48	1.9574	0.0109	0.0109	0.5054	1.0217
3	36	8.2286	0.0035	0.0035	0.5932	1.4583	7	42	5.2683	0.0046	0.0046	0.5325	1.1389							
3	39	8.8947	0.003	0.003	0.5938	1.4615	8	16	1.8667	0.0357	0.0357	0.5172	1.0714							

**Table 2: Optimal fat-tree designs  $FT([n_1, n_2])$  for skewness  $\sigma = 2$ .**

$n_1$	$n_2$	$c_1/c_2$	$s_1$	$s_2$	$\eta$	$\tau$	$n_1$	$n_2$	$c_1/c_2$	$s_1$	$s_2$	$\eta$	$\tau$	$n_1$	$n_2$	$c_1/c_2$	$s_1$	$s_2$	$\eta$	$\tau$
2	4	1.0	0.25	0.5	0.6667	2.0	3	42	7.2593	0.0026	0.0051	0.6585	1.9286	8	24	2.52	0.0159	0.0317	0.5435	1.1905
2	6	1.2857	0.1111	0.2222	0.7	2.3333	3	45	7.7586	0.0022	0.0044	0.6591	1.9333	8	32	3.2941	0.0089	0.0179	0.5484	1.2143
2	8	1.6	0.0625	0.125	0.7143	2.5	3	48	8.2581	0.002	0.0039	0.6596	1.9375	8	40	4.0698	0.0057	0.0114	0.5513	1.2286
2	10	1.9231	0.04	0.08	0.7222	2.6	4	8	1.5	0.0833	0.1667	0.5714	1.3333	8	48	4.8462	0.004	0.0079	0.5532	1.2381
2	12	2.25	0.0278	0.0556	0.7273	2.6667	4	12	2.0769	0.037	0.0741	0.5909	1.4444	9	18	1.7778	0.0312	0.0625	0.5294	1.125
2	14	2.5789	0.0204	0.0408	0.7308	2.7143	4	16	2.6667	0.0208	0.0417	0.6	1.5	9	27	2.5714	0.0139	0.0278	0.5385	1.1667
2	16	2.9091	0.0156	0.0312	0.7333	2.75	4	20	3.2609	0.0133	0.0267	0.6053	1.5333	9	36	3.3684	0.0078	0.0156	0.5429	1.1875
2	18	3.24	0.0123	0.0247	0.7353	2.7778	4	24	3.8571	0.0093	0.0185	0.6087	1.5556	9	45	4.1667	0.005	0.01	0.5455	1.2
2	20	3.5714	0.01	0.02	0.7368	2.8	4	28	4.4545	0.0068	0.0136	0.6111	1.5714	10	20	1.8	0.0278	0.0556	0.5263	1.1111
2	22	3.9032	0.0083	0.0165	0.7381	2.8182	4	32	5.0526	0.0052	0.0104	0.6129	1.5833	10	30	2.6129	0.0123	0.0247	0.5345	1.1481
2	24	4.2353	0.0069	0.0139	0.7391	2.8333	4	36	5.6512	0.0041	0.0082	0.6143	1.5926	10	40	3.4286	0.0069	0.0139	0.5385	1.1667
2	26	4.5676	0.0059	0.0118	0.74	2.8462	4	40	6.25	0.0033	0.0067	0.6154	1.6	11	22	1.8182	0.025	0.05	0.5238	1.1
2	28	4.9	0.0051	0.0102	0.7407	2.8571	4	44	6.8491	0.0028	0.0055	0.6163	1.6061	11	33	2.6471	0.0111	0.0222	0.5312	1.1333
2	30	5.2326	0.0044	0.0089	0.7414	2.8667	4	48	7.4483	0.0023	0.0046	0.617	1.6111	11	44	3.4783	0.0063	0.0125	0.5349	1.15
2	32	5.5652	0.0039	0.0078	0.7419	2.875	5	10	1.6	0.0625	0.125	0.5556	1.25	12	24	1.8333	0.0227	0.0455	0.5217	1.0909
2	34	5.898	0.0035	0.0069	0.7424	2.8824	5	15	2.25	0.0278	0.0556	0.5714	1.3333	12	36	2.6757	0.0101	0.0202	0.5286	1.1212
2	36	6.2308	0.0031	0.0062	0.7429	2.8889	5	20	2.9091	0.0156	0.0312	0.5789	1.375	12	48	3.52	0.0057	0.0114	0.5319	1.1364
2	38	6.5636	0.0028	0.0055	0.7432	2.8947	5	25	3.5714	0.01	0.02	0.5833	1.4	13	26	1.8462	0.0208	0.0417	0.52	1.0833
2	40	6.8966	0.0025	0.005	0.7436	2.9	5	30	4.2353	0.0069	0.0139	0.5862	1.4167	13	39	2.7	0.0093	0.0185	0.5263	1.1111
2	42	7.2295	0.0023	0.0045	0.7439	2.9048	5	35	4.9	0.0051	0.0102	0.5882	1.4286	14	28	1.8571	0.0192	0.0385	0.5185	1.0769
2	44	7.5625	0.0021	0.0041	0.7442	2.9091	5	40	5.5652	0.0039	0.0078	0.5897	1.4375	14	42	2.7209	0.0085	0.0171	0.5244	1.1026
2	46	7.8955	0.0019	0.0038	0.7444	2.913	5	45	6.2308	0.0031	0.0062	0.5909	1.4444	15	30	1.8667	0.0179	0.0357	0.5172	1.0714
2	48	8.2286	0.0017	0.0035	0.7447	2.9167	6	12	1.6667	0.05	0.1	0.5455	1.2	15	45	2.7391	0.0079	0.0159	0.5227	1.0952
3	6	1.3333	0.125	0.25	0.6	1.5	6	18	2.3684	0.0222	0.0444	0.5588	1.2667	16	32	1.875	0.0167	0.0333	0.5161	1.0667
3	9	1.8	0.0556	0.1111	0.625	1.6667	6	24	3.0769	0.0125	0.025	0.5652	1.3	16	48	2.7551	0.0074	0.0148	0.5213	1.0889
3	12	2.2857	0.0312	0.0625	0.6364	1.75	6	30	3.7879	0.008	0.016	0.569	1.32	17	34	1.8824	0.0156	0.0312	0.5152	1.0625
3	15	2.7778	0.02	0.04	0.6429	1.8	6	36	4.5	0.0056	0.0111	0.5714	1.3333	18	36	1.8889	0.0147	0.0294	0.5143	1.0588
3	18	3.2727	0.0139	0.0278	0.6471	1.8333	6	42	5.2128	0.0041	0.0082	0.5732	1.3429	19	38	1.8947	0.0139	0.0278	0.5135	1.0556
3	21	3.7692	0.0102	0.0204	0.65	1.8571	6	48	5.9259	0.0031	0.0063	0.5745	1.35	20	40	1.9	0.0132	0.0263	0.5128	1.0526
3	24	4.2667	0.0078	0.0156	0.6522	1.875	7	14	1.7143	0.0417	0.0833	0.5385	1.1667	21	42	1.9048	0.0125	0.025	0.5122	1.05
3	27	4.7647	0.0062	0.0123	0.6538	1.8889	7	21	2.4545	0.0185	0.037	0.55	1.2222	22	44	1.9091	0.0119	0.0238	0.5116	1.0476
3	30	5.2632	0.005	0.01	0.6552	1.9	7	28	3.2	0.0104	0.0208	0.5556	1.25	23	46	1.913	0.0114	0.0227	0.5111	1.0455
3	33	5.7619	0.0041	0.0083	0.6563	1.9091	7	35	3.9474	0.0067	0.0133	0.5588	1.2667	24	48	1.9167	0.0109	0.0217	0.5106	1.0435
3	36	6.2609	0.0035	0.0069	0.6571	1.9167	7	42	4.6957	0.0046	0.0093	0.561	1.2778							
3	39	6.76	0.003	0.0059	0.6579	1.9231	8	16	1.75	0.0357	0.0714	0.5333	1.1429							