

# Flattened Clos: Designing High-performance Deadlock-free Expander Data Center Networks Using Graph Contraction

## Abstract

Clos networks have witnessed the successful deployment of RoCE in production data centers. However, as DCN bandwidth keeps increasing, building Clos networks is becoming cost-prohibitive and thus the more cost-efficient expander graph has received much attention in recent literature. Unfortunately, the existing expander graphs' topology and routing designs may contain Cyclic Buffer Dependency (CBD) and incur deadlocks in PFC-enabled RoCE networks.

We propose Flattened Clos (FC), a topology/routing co-designed approach, to eliminate the PFC-induced deadlocks in expander networks. FC's topology and routing are designed in three steps: 1) logically divide each ToR switch into  $k$  virtual layers and establish connections only between adjacent virtual layers; 2) generate *virtual up-down* paths for routing; 3) flatten the virtual multi-layered network and the virtual up-down paths using graph contraction. We rigorously prove that FC's design is deadlock-free and validate this property using a real testbed and packet-level simulation. Compared to expander graphs with the edge-disjoint-spanning-tree (EDST) based routing (a state-of-art CBD-free routing algorithm for expander graphs), FC reduces the average hop count by at least 50% and improves network throughput by  $2 - 10 \times$  or more. Compared to Clos networks with up-down routing, FC increases network throughput by  $1.1 - 2 \times$  under all-to-all and uniform random traffic patterns.

## 1 Introduction

Driven by the need of low latency, high throughput and low CPU overhead, large Internet service providers such as Microsoft and Alibaba have deployed RDMA over Commodity Ethernet (RoCE) [10, 15] in their Clos data centers. RoCE requires a lossless network for optimal performance. To avoid packet loss in Ethernet, Priority-based Flow Control (PFC) is usually enabled to perform a hop-by-hop flow control to avoid exhausting switch buffers by upstreaming flows. However, enabling PFC introduces the risk of deadlocks, especially for

the large-scale deployment of RoCEv2. Thanks to the layered structure of Clos data centers, the up-down routing in Clos networks can prevent deadlocks with proper safety mechanisms under normal operations [15] and failure scenarios [19].

However, as the data center traffic and the network bandwidth keep increasing, building Clos topologies is becoming cost-prohibitive [3]. In order to reduce the network cost, flatter expander graphs, such as Jellyfish [38], SlimFly [4], Xpander [42], FatClique [45], etc., have been proposed to build data centers. A recent study [28] shows that a full throughput expander uses 25% fewer switches than a full throughput Clos. Note that the throughput values of expander graphs are attained using a multi-commodity flow formulation based on the K-Shortest Path (KSP) routing [44]. Unfortunately, the KSP routing in expander graphs is not CBD-free, and thus could incur severe PFC deadlocks. Therefore, the performance-gain or cost-reduction of expander graphs over Clos becomes questionable for RoCEv2 traffic.

The key to supporting RoCE in expander graphs is to eliminate Cyclic Buffer Dependency (CBD). Approaches to eliminate CBD can be generally grouped into three classes. The first approach is to assign different lossless priorities for packets at different hops [9, 11, 22]. This approach has been widely adopted in HPCs, in which the underlying Infiniband network supports 15 lossless priorities (*a.k.a.* Virtual Channel). However, due to the limited switch buffer space, data center switches can support at most two or three lossless priorities [15]. The second approach is to disable PFC and redesign RoCE to work with lossy networks, e.g., NDP [16], IRN [27], FatPaths [5], etc. However, lossy RoCE requires hardware support. For example, Mellanox ConnectX-4 onwards NICs support lossy RoCE, but Mellanox ConnectX-3 NICs do not. In addition, lossy RoCE may incur higher latency for mice flows, especially when a sender has to rely on a timeout to retransmit a lost packet. iWarp [33] is another RDMA technology that runs on lossy networks. However, its performance is poor because it relies on TCP to guarantee lossless delivery.

The third approach is to design a routing solution that is fundamentally free of CBD, just as the up-down routing in

Clos. Along this direction, TCP Bolt [39, 40] was proposed, and the key idea is to find as many edge disjoint spanning trees (EDST) in an expander graph as possible and then route each packet in one of the spanning tree from its source to its destination. The EDST-based routing is CBD-free, but its throughput performance is poor. The key reason is that the EDST-based routing cannot effectively utilize all the network resources: 1) the average path length is usually large and increases quickly with network size; 2) some network links could remain idle as they do not belong to any EDST.

Our work called **Flattened Clos (FC)** offers a novel topology-routing co-design to eliminate CBDs in RoCE networks. FC’s topology is essentially a random regular graph that is mappable to a multi-layered topology. We construct FC’s topology in two steps: 1) virtually split each ToR switch into  $k$  virtual switches, each of which belongs to a virtual layer, and 2) randomly interconnect the layer- $i$  ( $i = 2, \dots, k - 1$ ) virtual switches to the layer- $(i - 1)$  and the layer- $(i + 1)$  virtual switches. The multi-layered virtual structure of FC allows performing up-down routing based on virtual layers. To this end, we propose the Edge-Disjoint Virtual Up-Down Routing for FC. For every source-destination pair, FC’s routing transforms the path-finding problem into a min-cost-flow problem and then finds the maximum number of edge-disjoint paths. We analyze FC’s design in the following aspects to demonstrate its feasibility:

1. We offer a theoretical guidance for choosing the right number  $k$  of virtual layers when constructing FC’s topology (see the strategy (\*) in Section 3.2.3), and validate the strategy via numerical analysis.
2. We prove that FC’s routing is CBD-free, and thus is deadlock-free. In fact, FC’s topology and routing paths can be viewed as contracted graphs of a virtual multi-layered network and virtual up-down paths. This *graph contraction* operation preserves the CBD-free property.
3. We show that FC’s cabling complexity can be dramatically reduced by introducing a layer of Patch Panels (PP) or Optical Circuit Switches (OCS) to interconnect all the ToR switches.
4. We demonstrate that FC outperforms expander graphs with EDST routing [39, 40] (the state-of-art CBD-free routing for expanders). Specifically, FC reduces the average hop count (AHC) by at least 50% and increases network throughput by  $2 - 10\times$  or more.
5. We compare the throughput performance between FC and cost-equivalent Clos networks with up-down routing. FC achieves  $1.1 - 2\times$  throughput for all-to-all and uniform random traffic patterns, but its near-worst throughput is lower. We argue that when OCSs are used to construct FC, vendors do not have to worry much about FC’s near-worst throughput. By simply generating a different

FC’s topology, one can avoid matching an FC’s topology with its near-worst traffic patterns.

6. We validate that FC is deadlock free using a small test bed and a packet-level simulator, even under extreme (but practical) cases where congestion control is disabled and switches are misconfigured with a very small PFC PAUSE threshold. In contrast, we see deadlocks triggered under shortest-path routing and thus ECMP&KSP routing is not safe.

## 2 Background & Motivation

### 2.1 Deploy RDMA over Ethernet in Clos

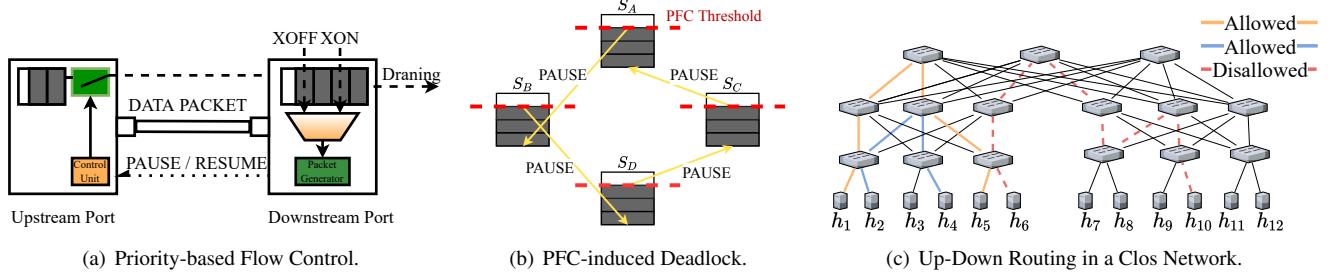
Clos network, a.k.a. fat-tree, was proposed for data center network (DCN) architecture in [2], and has become the de-facto standard for large service providers, such as Google [37], Microsoft [14], Facebook [36], etc. TCP/IP is the dominant transport/network stack in today’s data centers. However, the traditional TCP/IP stack cannot offer high throughput ( $> 40\text{Gbps}$  or more) and ultra-low latency ( $< 10\text{us}$  per hop) for modern data center applications such as cloud storage, deep learning framework and database [15, 24, 48]. Therefore, data center operators, e.g., Microsoft [15], Alibaba [24], etc., have started large-scale deployment of RDMA in Clos data centers to attain better network performance.

RDMA is a hardware offloading technology that offers several benefits such as high throughput, low latency and low CPU overhead by bypassing the host networking stack. HPC community has long used RDMA in special-purpose clusters, and deployed RDMA using Infiniband (IB) technology. However, modern data centers are built with IP and Ethernet technologies. For technical and economical reasons, RoCE was proposed for RDMA deployment in data centers.

The commonly used RoCE protocol is RoCEv2. RoCEv2 encapsulates an RDMA transport packet within a UDP packet to be compatible with the existing networking infrastructure of data centers. RoCEv2 was initially designed to run on a lossless network, which can be guaranteed by enabling the Priority-based Flow Control (PFC) [20]. Admittedly, there have been advanced RoCE designs, e.g., Resilient RoCE, IRN [27], etc., that could work with a lossy network. However, supporting RoCE in lossy networks requires handling packet retransmission using time out, selective ack, etc., which may not only complicate the NIC design, but also hurt network latency and throughput performance. Therefore, we still focus on lossless networks to support RoCEv2 in this paper.

#### 2.1.1 Priority-based Flow Control (PFC)

PFC is a hop-by-hop flow control approach to prevent switch buffer overflow, which is the primary cause of packet loss in data centers. As shown in Fig. 1(a), the downstream switch



**Figure 1:** Key Technologies of Supporting RoCE in a lossless Clos network.

sends a PAUSE frame to its upstream switch when its ingress queue length exceeds a certain threshold (**XOFF**). The upstream switch stops transmission after receiving the PAUSE frame. A RESUME frame is sent when the downstream queue drains below another threshold (**XON**). It takes some time for the upstream switch to react to the PAUSE frame and stop transmission. So the downstream switch needs to reserve some buffer space to accommodate the packets sent by the upstream switch during this time. The buffer space is called *headroom*. PFC can guarantee zero packet loss when the headroom size is configured correctly. Typically, data center switches can support at most two or three lossless priorities [15] due to the buffer size limit.

### 2.1.2 PFC-induced Deadlocks

PFC can raise some performance issues such as unfairness, PFC storms and deadlocks [10, 15, 24, 48]. Specifically, the PFC-induced deadlocks may hinder the large-scale deployment of RoCEv2. When cyclic buffer dependency (CBD) exists, deadlocks can be triggered by PFC PAUSES [18], causing packets to wait indefinitely for buffer resources [40]. As shown in Fig. 1(b), four switches  $S_A, S_B, S_C, S_D$  have reached the PFC threshold and start to send PAUSE frames; then the network is trapped into a deadlock and no switch can make any progress. Note that, the PFC-induced deadlock cannot go away once occurs even if we restart all the servers.

### 2.1.3 Avoiding Deadlocks in Clos Networks

Large vendors have gained years of experience in deploying RDMA in Clos data centers [15]. The following strategies are adopted to avoid deadlocks:

1. Perform up-down routing, which is CBD and deadlock-free under normal network conditions in Clos networks. (Note that containing a CBD is a necessary condition to have deadlocks.) As shown in Fig. 1(c), the paths of  $h_1 \rightarrow h_5$  and  $h_2 \rightarrow h_4$  obey the “up-down” rule and are allowed; but the path of  $h_6 \rightarrow h_{10}$  contains a “down-up” segment and thus is not allowed.

2. Do not put multicast and broadcast packets into lossless classes. It was reported in [15] that ARP broadcasts+up-down routing can cause PFC deadlocks.
3. Use a different lossless class for rerouted packets upon network failures. [19] shows that packet rerouting may break the “up-down” rule and trigger PFC deadlocks.

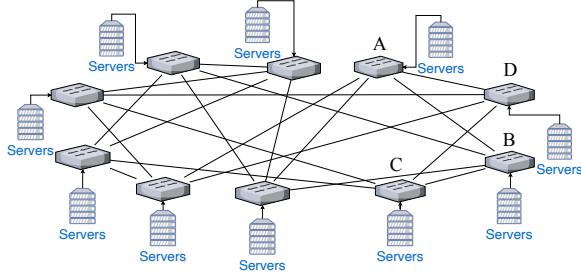
## 2.2 From Clos to Expander

Despite of the success of deploying Clos data centers, however, a Clos network is inherently suboptimal in terms of bandwidth provision. As the Ethernet speed keeps increasing, the network cost, especially the power consumption of Clos networks, is becoming prohibitively high [3]. To reduce the network cost, researchers have started seeking for more cost-effective network architectures.

One of the promising alternative for DCNs is expander graph. As shown in Fig. 2, expander graphs adopt a flat topology design: servers connect to ToR switches and these ToRs are directly interconnected without a layered structure. Examples of expander graphs include Jellyfish [38], SlimFly [4], Xpander [42], FatClique [45], etc. Expander graphs is more cost-effective for bandwidth provision than Clos networks. Using KSP routing, a full throughput expander uses 25% fewer switches than a full throughput Clos [28]. The network performance of expander graphs was also studied under other routing protocols, including ECMP, Valient Load Balancing (VLB) and a hybrid of the two [23]. However, none of these widely studied routing strategies is CBD-free.

### 2.2.1 ECMP/KSP are not CBD-free in Expanders

Consider the expander graph in Fig. 2. This expander is a random regular graph with 4 inter-ToR links per ToR. Consider the four ToRs  $A, B, C, D$  and the shortest paths for  $A \rightarrow C$ ,  $B \rightarrow D$ ,  $C \rightarrow A$ ,  $D \rightarrow B$ . Assume that there is a flow routed along the shortest path  $A \rightarrow B \rightarrow C$ . Then, if the egress port at link  $(B, C)$  is paused, the egress port at link  $(A, B)$  will be paused. If there is another flow routed along the shortest path  $D \rightarrow A \rightarrow B$ , since the egress port at link  $(A, B)$  is



**Figure 2:** An expander graph.

paused, the egress port at link  $(D, A)$  will also be paused. Similarly, if we have another two flows routed along the shortest paths  $C \rightarrow D \rightarrow A$  and  $B \rightarrow C \rightarrow D$ , then the egress ports at link  $(C, D)$  and link  $(B, C)$  will be paused. Now, we find a CBD in this expander graph under shortest-path routing. To sum up, when there are 4 flows routed along the paths  $A \rightarrow B \rightarrow C, D \rightarrow A \rightarrow B, C \rightarrow D \rightarrow A$  and  $B \rightarrow C \rightarrow D$ , if one of the egress ports  $(A, B), (B, C), (C, D), (D, A)$  is paused for a sufficiently long time, a deadlock will be triggered.

The above analysis indicates that shortest-path routing is not CBD-free. Now we consider ECMP and KSP routings. ECMP uniformly split traffic among all the shortest paths, while KSP split traffic among the first  $K$  shortest paths. (To improve network performance under ECMP/KSP, one can also optimize the path weights using a multi-commodity flow formulation.) Under ECMP or KSP routing, it is still possible to have four flows taking the paths  $A \rightarrow B \rightarrow C, D \rightarrow A \rightarrow B, C \rightarrow D \rightarrow A$  and  $B \rightarrow C \rightarrow D$  in the above example. Therefore, both ECMP and KSP routings are not CBD-free. Using the same approach, we can prove that the VLB routing and the hybrid of ECMP&VLB in [23] are not CBD-free, either.

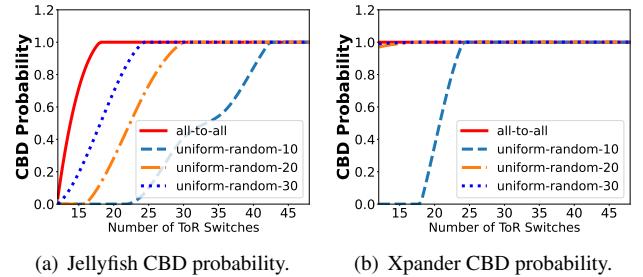
### 2.2.1.1 Probability of Containing CBDs

We further analyze the probability of an expander graph containing CBDs under different traffic patterns. We generate two classes of expander graphs, Jellyfish [38] and Xpander [42]. In each expander graph, each ToR switch has 5 ports connected to other ToRs. For each expander graph, we evaluate two classes of traffic patterns under shortest-path routing (the algorithm that determines if a set of paths is CBD free in a given topology is offered in Appendix A.2):

**All to All:** Every source-destination pair has an on-going flow. This represents the most-likely case of having CBDs.

**Uniform Random- $p$ :** Every ToR randomly picks  $p$  fraction of ToRs to communicate. This represents practical DCN traffic patterns in which the majority of traffic of a server is often destined to a few racks [36].

The results are depicted in Fig. 3. We can see that as the number of ToRs increases, the CBD probability quickly increases to one and Xpander graphs are more likely to encounter CBDs than Jellyfish graphs (random regular graphs). Note that even under shortest-path routing (ECMP), the CBD



(a) Jellyfish CBD probability.

(b) Xpander CBD probability.

**Figure 3:** Jellyfish and Xpander CBD analysis.

probability becomes one with only tens of ToR switches. The risk of deadlocks under other routing algorithms, including KSP, VLB, etc., would be even higher.

## 3 Flattened Clos

We propose a new class of expander graphs, called **Flattened Clos (FC)**, for efficient deadlock prevention. Our design is motivated by the CBD-free up-down routing in the multi-layered networks. FC is a topology built on top of ToR switches. The key idea of FC is to split each ToR switch into a few virtual switches, and assign each virtual switch a virtual layer id. By creating links only between virtual switches in adjacent virtual layers, FC can adopt virtual up-down routing to avoid deadlocks. The detailed topology and routing designs of FC are described below.

### 3.1 Topology

We study a data center network with  $N$  ToR (Top of Rack) switches  $S = \{S_1, S_2, \dots, S_N\}$ . Each switch has  $p = s + h$  ports, with  $p$  of which connected to the hosts and  $s$  of which connected to other ToR switches. We construct a Flattened Clos topology in two steps:

**Step 1: Splitting Virtual Switches.** To create an FC with  $k$  virtual layers, we logically split each switch  $S_i, i = 1, 2, \dots, N$  into  $k$  virtual switches, and label these virtual switches as  $S_i^1, S_i^2, \dots, S_i^k$ . The virtual switch  $S_i^j$  belongs to the  $j$ -th layer, and has  $l_j$  number of links to connect to other ToR switches. The total number of links of the virtual switches  $S_i^1, S_i^2, \dots, S_i^k$  is equal to the total number of links of  $S_i$  that connect to other switches, i.e.,

$$\sum_{j=1}^k l_j = s. \quad (1)$$

**Step 2: Random Wiring.** For each  $j = 1, 2, \dots, k - 1$ , we randomly generate a bipartite graph between the virtual switches in layer  $j$  and the virtual switches in layer  $j + 1$ . Let  $a_j, j = 1, 2, \dots, k - 1$  be the degree of each virtual switch in the  $j$ -th random bipartite graph. We must have

$$l_1 = a_1, l_2 = a_1 + a_2, \dots, l_{k-1} = a_{k-2} + a_{k-1}, l_k = a_{k-1}. \quad (2)$$

When we generate random bipartite graphs, we never create links between  $S_i^j$  and  $S_i^{j+1}$  for  $i = 1, 2, \dots, N, j = 1, 2, \dots, k - 1$ . The reason is that  $S_i^j$  and  $S_i^{j+1}$  actually belong to the same switch, and there is no need to create a link in between.

### 3.1.1 Theoretical Topology Properties of FC

In an FC's topology, each ToR switch has  $s$  links connected to other ToRs. Thus, FC falls into the category of random  $s$ -regular graphs. Here, we restate some useful theoretical results for random regular graphs (RRG) in literature, which applies to FC, as well.

We represent a network by  $G = (V, E)$ , where  $V$  is the vertex set and  $E$  is the edge set. The bisection bandwidth of  $G$  can be characterized by *Edge Expansion*, which is defined as

$$EE(G) = \min_{|S| \leq \frac{N}{2}} \frac{|\partial S|}{|S|},$$

where  $N$  is the number of vertices in  $V$ ,  $S$  is a subset of  $V$ ,  $|S|$  is the size of  $S$ ,  $\partial S$  is the set of edges leaving  $S$ . The Edge Expansion of an  $s$ -regular graph is upper bounded by  $s/2$  [42]. The following theorem indicates that random regular graphs attain near-optimal edge expansion.

**Theorem 1** (Near-optimal Edge Expansion [6]) *For every  $s \geq 3$  and  $0 < \eta < 1$  satisfying*

$$2^{4/s} < (1 - \eta)^{1-\eta} (1 + \eta)^{1+\eta},$$

*almost every  $s$ -regular graph  $G$  has its edge expansion*

$$EE(G) \geq (1 - \eta)s/2.$$

Given a traffic matrix  $T = [t_{uv}]$ , where  $t_{uv}$  is the amount of requested flows from ToR switch  $u$  to ToR switch  $v$ . The throughput  $\alpha(G, T)$  of a network  $G$  under the traffic matrix  $T$  is defined as the maximum value  $\theta(T)$  for which  $T \cdot \theta(T)$  is feasible in  $G$ . The following two theorems guarantee that random regular graphs achieve good throughput under both uniform and adversarial patterns.

**Theorem 2** (High throughput under all-to-all pattern [42]): *For the all-to-all traffic pattern  $T_{\text{all-to-all}}$ , almost every  $s$ -regular graph  $G$  achieves a throughput*

$$\alpha(G, T_{\text{all-to-all}}) \geq \frac{1}{O(\log s)} \alpha(G^*, T_{\text{all-to-all}}),$$

*where  $G^*$  is the  $s$ -regular graph that attains the optimal throughput under  $T_{\text{all-to-all}}$ .*

**Theorem 3** (Resilience to adversarial patterns [42]): *For almost every  $s$ -regular graph  $G$  and every traffic pattern  $T$ , the throughput  $\alpha(G, T) \geq \frac{1}{O(\log N)} \alpha(G^*, T)$ , where  $G^*$  is the  $s$ -regular graph that attains the optimal throughput under  $T$ .*

## 3.2 Routing

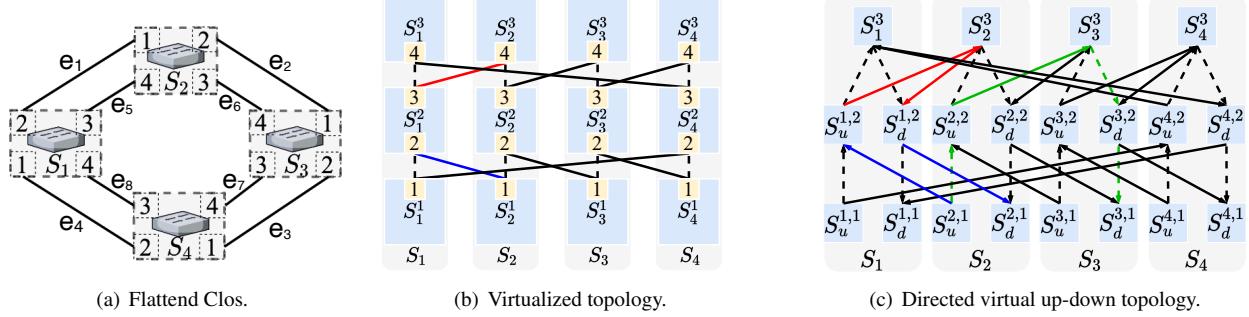
### 3.2.1 Edge-disjoint Virtual Up-down Routing

Although FC's topology exhibits high network throughput in theory, such a throughput may not be achievable in PFC-enabled RoCE networks due to the potential risk of deadlocks. To completely eliminate the risk of deadlocks, we propose the CBD-free **Edge-disjoint Virtual Up-down Routing**. This routing strategy computes paths in three steps:

**Step 1: Construct a Multi-layered Virtual Topology.** According to the construction of FC's topology, each FC's topology is mappable to a multi-layered topology. Consider the toy example in Fig. 4(a). While we construct this topology, we have virtually divided each ToR switch  $S_i$  into  $k = 3$  sub-switches  $S_i^1, S_i^2$  and  $S_i^3$ . The first port of  $S_i$  belongs to  $S_i^1$ ; the second and the third ports belong to  $S_i^2$ ; the fourth port belongs to  $S_i^3$ . It is easy to check that each edge in Fig. 4(a) corresponds to a solid line in Fig. 4(b). Note that the  $k$  sub-switches  $S_i^1, S_i^2, \dots, S_i^k$  can communicate with each other, because they belong to the same physical switch. Hence, we also create an edge between  $S_i^j$  and  $S_i^{j+1}$  for every  $j = 1, 2, \dots, k - 1$  (see the dashed lines) in Fig. 4(b).

**Step 2: Construct a Directed Virtual Up-Down Topology.** Our objective is to find the maximum number of virtual up-down paths in the multi-layered virtual topology. To enforce this “up-down” constraint, we further convert the undirected multi-layered graph in Fig. 4(b) to a directed virtual up-down graph in Fig. 4(c). Specifically, we first split each virtual node  $S_i^j$  ( $j < k$ ) into one “up node”  $S_u^{i,j}$  and one “down node”  $S_d^{i,j}$  in the directed up-down graph. (Note that we do not split the top layer virtual nodes  $S_i^k$ .) We then map each link in Fig. 4(b) to two directed links in Fig. 4(c): each undirected link  $(S_{i_1}^{k-1}, S_{i_2}^k)$  in the top layer (see the red line in Fig. 4(b) as an example) is mapped to two directed links  $(S_u^{i_1,k-1}, S_{i_2}^k)$  and  $(S_{i_2}^k, S_d^{i_1,k-1})$ ; each undirected link  $(S_{i_1}^{j-1}, S_{i_2}^j)$  ( $j = 2, \dots, k - 1$ ) (see the blue line in Fig. 4(b) as an example) is mapped to two directed links  $(S_u^{i_1,j-1}, S_{i_2}^{j-1})$  and  $(S_{i_2}^{j-1}, S_d^{i_1,j-1})$ .

**Step 3: Compute CBD-free Paths.** For every source-destination pair  $(S_i, S_j)$ , we first find a path set  $\mathcal{P}_{ij}$  with the maximum number of virtual up-down paths from the node  $S_u^{i,1}$  to the node  $S_d^{j,1}$  in the directed virtual up-down topology using min-cost max-flow (see Appendix A.1 for more details). In this set  $\mathcal{P}_{ij}$  of paths, each solid link is used at most once while each dashed link can be used multiple times. Then, for every path  $P \in \mathcal{P}_{ij}$ , we map it to a path in FC's topology (Fig. 4(a)). For example, as shown in Fig. 4(c), we find one up-down path  $S_u^{2,1} \rightarrow S_u^{2,2} \rightarrow S_3^3 \rightarrow S_d^{3,2} \rightarrow S_d^{3,1}$  (marked with green) for the source-destination pair  $(S_1, S_2)$ . Since  $S_u^{2,1}, S_u^{2,2}$  are from the ToR switch  $S_2$  and  $S_3^3, S_d^{3,2}, S_d^{3,1}$  are from the ToR switch  $S_3$ , this path can be contracted to  $S_2 \rightarrow S_3$  in the FC's topology. Since each solid link is used at most once in  $\mathcal{P}_{ij}$ , the resulting paths in the FC's topology must be edge-disjoint.



**Figure 4:** FC’s topology and routing design.

### 3.2.2 FC’s Routing is CBD Free

In FC’s edge-disjoint virtual up-down routing, we first compute an up-down path set  $\mathcal{P}_{ij}$  based on the directed virtual up-down topology, and then contract all the paths in  $\mathcal{P}_{ij}$  to obtain the final paths for FC’s topology. Let  $\mathcal{P} = \cup_{i,j} \mathcal{P}_{ij}$  be the set of virtual up-down paths obtained from the directed virtual up-down topology. According to Theorem 8 in Appendix A.2,  $\mathcal{P}$  is CBD free. In order to prove that the final set of paths in FC’s topology is CBD free, we need the following definition and lemma (see Appendix A.2.1 for the proof).

**Definition 1** Given a set of nodes  $V$ ,  $\{V_1, V_2, \dots, V_m\}$  is called a partition of  $V$ , if the following conditions are met: 1)  $V_{m_1} \cap V_{m_2} = \emptyset$  for every  $m_1 \neq m_2$ ; 2)  $V_1 \cup V_2 \cup \dots \cup V_m = V$ .

**Lemma 4** Given a graph  $G(V, E)$ , a path set  $\mathcal{P}$  and a partition  $\{V_1, V_2, \dots, V_m\}$  of  $V$ , a graph and path set pair  $(\hat{G}(\hat{V}, \hat{E}), \hat{\mathcal{P}})$  is called a contraction of  $(G(V, E), \mathcal{P})$  if

1. every node in  $\hat{V}_i \in \hat{V}$  corresponds to the vertex set  $V_i$ ;
2. the number of edges between  $\hat{V}_i$  and  $\hat{V}_j$  is the same as the total number of edges between  $V_i$  and  $V_j$  in  $G(V, E)$ ;
3. each path  $\hat{P} \in \hat{\mathcal{P}}$  is a contraction of a path  $P \in \mathcal{P}$ , i.e.,  $\hat{P}$  is obtained by first replacing each vertex in  $P$  by a vertex in  $\hat{V}$  and then removing cycles and duplicated vertices.

Then, if the path set  $\mathcal{P}$  is CBD-free in  $G(V, E)$ , the path set  $\hat{\mathcal{P}}$  must be CBD-free in  $\hat{G}(\hat{V}, \hat{E})$ .

Apparently, FC’s topology and routing path set can be viewed as a contraction of the directed virtual up-down topology and the corresponding virtual up-down path set  $\mathcal{P}$ . Since the path set  $\mathcal{P}$  is CBD free in the directed virtual up-down topology, then according to Lemma 4, we immediately know that FC’s routing path set is CBD free.

### 3.2.3 How Routing Affects FC’s Topology Design?

We have described FC’s routing and topology designs. Note that there is a critical parameter  $k$  in the design. If  $k$  is not

properly chosen, FC’s routing policy may not be able to find a path for some switch pair, thus hurting the connectivity of FC. In this section, we offer a theoretical guideline to determine the number of virtual layers in FC.

**Lemma 5** Let  $x$  be the number of ancestors in the virtual layer  $k$  for each layer-1 virtual node. If  $x > \sqrt{(2+\epsilon)N \ln N}$ , where  $\epsilon > 0$  is an infinitesimal value, then as  $N \rightarrow +\infty$ , with probability 1, every pair of layer-1 virtual nodes has a common ancestor in the virtual layer  $k$ .

**Proof 1** We use  $A_{ij}$  to denote the event that the virtual nodes  $S_i^1$  and  $S_j^1$  have no common ancestor in the virtual layer  $k$ . Then, the probability that  $A_{ij}$  happens is

$$\begin{aligned} P(A_{ij}) &= \frac{C_{N-x}^x}{C_N^x} \leq \left(1 - \frac{x}{N}\right)^x \\ &< \left(1 - \frac{\sqrt{(2+\epsilon)\ln N}}{\sqrt{N}}\right)^{\sqrt{(2+\epsilon)N \ln N}} \\ &= \left(\left(1 - \frac{\sqrt{(2+\epsilon)\ln N}}{\sqrt{N}}\right)^{\frac{\sqrt{N}}{\sqrt{(2+\epsilon)\ln N}}}\right)^{(2+\epsilon)\ln N} \\ &< (1/e)^{(2+\epsilon)\ln N} = N^{-(2+\epsilon)}. \end{aligned}$$

Let  $A$  be the event that at least one pair of virtual nodes in layer 1 has no common ancestor in layer  $k$ . Then,

$$\begin{aligned} P(A) &= P(\cup_{i \neq j} A_{ij}) \leq \sum_{i \neq j} P(A_{ij}) \\ &= \frac{N(N-1)}{2} \times N^{-(2+\epsilon)} < \frac{1}{2} N^{-\epsilon}. \end{aligned}$$

Then,  $\lim_{N \rightarrow +\infty} P(A) = 0$ . This completes the proof.

Based on FC’s routing, it is easy to calculate that the number of distinct up-paths from a virtual node in layer 1 is  $(a_1 + 1)(a_2 + 1) \cdots (a_{k-1} + 1)$ , which is an upper bound of the number of ancestors in layer  $k$ . According to Lemma 5, we can choose a  $k$  such that

$$(a_1 + 1)(a_2 + 1) \cdots (a_{k-1} + 1) > \sqrt{(2+\epsilon)N \ln N}. \quad (3)$$

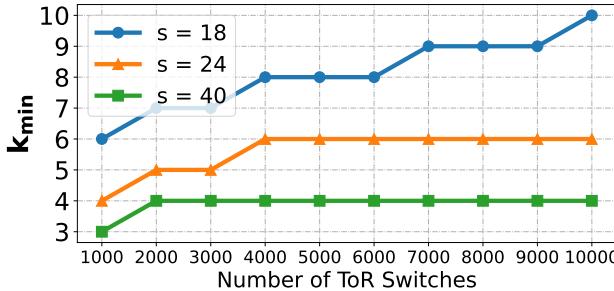
**Table 1:** Choosing the right  $k$  for FC.

Number of Switches	Number of Servers	$k_{\min}$	$k$	Average Number of Edge Disjoint Up-down Paths	Average Path Length of Edge Disjoint Up-down Paths	Minimum Number of Edge Disjoint Up-down Paths
500	12000	3	3	10.05	4.29	4.00
			4	16.08	4.57	12.00
1000	24000	3	3	7.10	4.43	1.00
			4	14.01	4.86	9.00
2000	48000	4	4	11.85	5.13	7.00
			5	15.77	5.36	11.00
			4	10.63	5.30	6.00
3000	72000	4	5	14.66	5.54	10.00
			4	9.17	5.52	3.00
5000	120000	4	5	13.28	5.75	9.00

According to Equation (1) and (2), it is easy to obtain  $\sum_{i=1}^{k-1} a_i = s/2$ . We could choose  $a_1, a_2, \dots, a_{k-1}$  to maximize the left hand side of (3), and obtain

$$\left(1 + \frac{s}{2(k-1)}\right)^{k-1} > \sqrt{(2+\varepsilon)N \ln N}. \quad (4)$$

Let  $k_{\min}$  be the smallest integer solution of (4). We could choose a  $k$  value around  $k_{\min}$ . As shown in Fig. 5,  $k_{\min}$  does not grow fast with respect to  $N$ .



**Figure 5:** Relationship between  $k_{\min}$  and network size  $N$ .

**Numerical Verification:** To verify the above theoretical result, we perform a numerical analysis using 64-port ToR switches. Each ToR switch has  $h = 24$  ports connected to the hosts and  $s = 40$  ports connected to other ToR switches. For different number of ToR switches ( $N = 500/1000/2000/3000/5000$ ), we choose  $k = k_{\min}, k_{\min} + 1$ , generate an FC’s topology and count the number of distinct virtual up-down paths. As shown in Table 1, as we increase  $k$ , more paths can be found for every source-destination ToR switch pairs. Note that the average path length increases with respect to  $k$ . Hence, it is better to choose a smaller  $k$ . On the other hand, if we choose  $k$  to be too small, some ToR switch pairs may not have sufficient number of distinct paths. Here we suggest a simple strategy that works well for FC:

**Strategy (\*):** Try  $k_{\min}$  first; if not working, try  $k_{\min} + 1$ .

For example, in the case where  $N = 1000, k = k_{\min} = 3$ , the minimum number of distinct paths between ToR pairs is 1.

This creates a bottleneck in the network. Hence,  $k = k_{\min} + 1 = 4$  will be chosen instead.

### 3.2.4 Computational Complexity of FC’s Routing

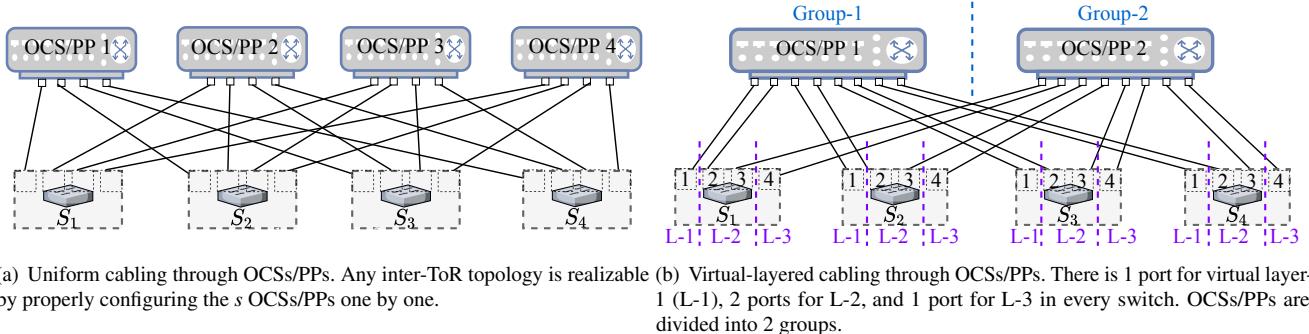
The main complexity comes from using the min-cost max-flow solver to find edge-disjoint virtual up-down paths. Given a graph  $G = (V, E)$  with  $n$  vertices and  $m$  edges, the computational complexity of the min-cost max-flow algorithm implemented in Ortools is  $O(m \cdot n^2 \cdot \log(n \cdot C))$ , where  $C$  is the value of the largest link cost in the graph [1] (in our case,  $C = 1$ ). If we choose the parameters  $k$  and  $a_1, a_2, \dots, a_{k-1}$  such that  $(a_1 + 1)(a_2 + 1) \cdots (a_{k-1} + 1) = \Theta(\sqrt{N \log N})$ , each virtual node in the first virtual layer will be able to reach at most  $\Theta(k\sqrt{N \log N})$  virtual nodes through at most  $\Theta(k\sqrt{N \log N})$  edges. When we compute edge disjoint paths from  $S_i$  to  $S_j$ , we only need to focus on a subgraph of the directed virtual up-down topology, which contains all the nodes reachable from  $S_u^{i,1}$  and  $S_d^{j,1}$ . This subgraph has  $\Theta(k\sqrt{N \log N})$  nodes and edges. Thus, the overall computational complexity is  $\Theta((k\sqrt{N \log N})^3 \cdot \log(k\sqrt{N \log N})) = \Theta(k^3 N^{3/2} (\log N)^{5/2})$ .

### 3.3 Cabling

FC adopts random wiring for its topology design. However, random wiring has long been criticized for its high cabling complexity [42, 45]. Indeed, if we directly connect different ToR switch pairs, the number of distinct fiber lengths would be in the order of  $\Theta(N^2)$ . Directly connecting ToR switches could also increase the management complexity when we perform data center expansion [47].

To reduce cabling complexity, motivated by TROD [7] and Google’s Jupiter data center [31], we propose to use a set of co-located optical circuit switches (OCS) or patch panels (PP) to interconnect different ToR pairs and form FC’s topology. Since these PPs/OCSs are co-located, the number of distinct fiber lengths reduces to  $\Theta(N)$ . Next, we offer two strategies to interconnect PPs/OCSs with ToR switches.

**Uniform Cabling (see Figure 6(a)):** Note that each ToR switch has  $s$  ports to be connected to other ToR switches. We use  $s$  OCSs/PPs, and construct a uniform bipartite graph



**Figure 6:** Example of cabling with the help of optical circuit switches (OCS) / patch panels (PP).

between ToR switches and OCSs/PPs. Under this cabling strategy, it was proven in [46] (see Lemma 4 and Theorem 5 therein) that any inter-ToR topology is realizable by properly configuring the  $s$  OCSs/PPs one by one. According to this fact, we could first generate an FC topology without considering the layer of OCSs/PPs, and then decompose this topology into  $s$  sub-topologies that can be mapped to each OCS/PP. This approach reduces cabling complexity. However, it encounters scalability challenge. Specifically, the port count of the commercially available OCSs/PPs is on the order of a few hundred. For example, a Calient s320 OCS [41] can offer 320 TX/RX ports. Thus, the number of ToR switches can be at most a few hundreds. Since each ToR switch typically connects to tens of servers, this uniform cabling strategy can support at most a few thousands of servers.

**Virtual-Layered Cabling (see Figure 6(b)):** Note that FC's topology is designed based on the concept of virtual layers. Assume that there are  $a_1$  ports for layer-1,  $a_1 + a_2$  ports for layer-2, ...,  $a_{k-1} + a_k$  ports for layer- $(k-1)$ , and  $a_k$  ports for layer- $k$ . We group all the OCSs/PPs into  $k-1$  groups, and connect  $2a_i$  ports of each ToR switch to the  $i$ -th OCS/PP group. In the  $i$ -th OCS/PP group, each OCS/PP have half of its ports connected to ToR switches' layer- $i$  ports and half of its ports connected to ToR switches' layer- $(i+1)$  ports. If we enforce that every OCS/PP should connect to all the ToR switches, we will encounter the same scalability challenge as the uniform cabling strategy. In the virtual-layered cabling strategy,  $2\eta a_i$  number of OCSs/PPs are used in the  $i$ -th group, and each ToR switch will randomly choose  $2a_i$  OCSs/PPs in group- $i$  to connect its layer- $i$  and layer- $(i+1)$  ports. Under this cabling strategy, the total number of ToR switches that can be supported becomes " $\eta \times$  port count of an OCS/PP". This strategy scales well. For example, if we use 320-port OCSs, 64-port ToR switches (assume that each ToR connects to 24 servers), and choose  $\eta = 20$ , then the maximum number of servers would be  $20 \times 320 \times 24 = 153600$ , which can definitely support a large-scale data center.

**Remark:** When  $\eta > 1$ , a cabling constraint is imposed to FC when we generate the topology between adjacent virtual

layers, i.e., not all topologies are realizable because the interconnection between ToRs and each group of OCSs/PPs is not uniform. Fortunately, as shown in Table 2, enabling this cabling constraint when generating FC's topology has little impact on FC's routing statistics.

## 4 Numerical Throughput Analysis

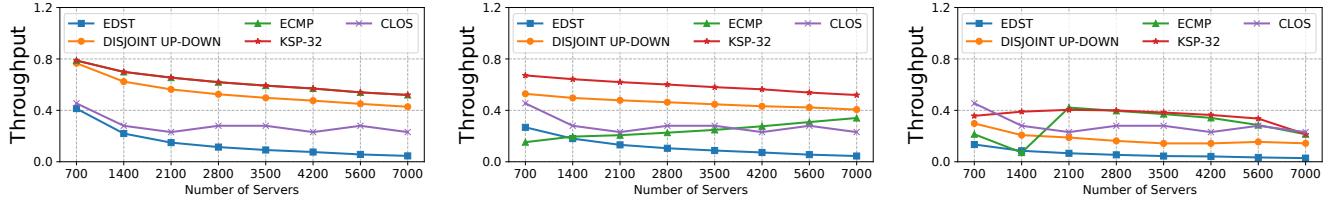
We numerically evaluate the throughput for FC in this section. We evaluate two scenarios. Due to space constraints, we only present one here and put the other one in Appendix A.4.1.

We generate FC's topologies of different sizes using up to 500 32-port ToR switches. Each ToR switch has 18 ports connected to other switches and 14 ports connected to servers. The number of virtual layers  $k$  is chosen based on the strategy (\*) in Section 3.2.3. For each FC's topology, we evaluate four routing strategies: 1) FC's edge-disjoint virtual up-down routing, 2) EDST routing, 3) ECMP or Shortest-Path routing, and 4) KSP routing. Given a traffic matrix  $T$ , we use a multi-commodity flow formulation to calculate the maximum throughput value  $\theta(T)$  such that  $T \cdot \theta(T)$  is feasible under the given topology and routing paths. (For ECMP, the throughput is also calculated based on the multi-commodity flow formulation. Evenly spreading traffic among all the shortest paths may yield very poor throughput.) In addition, for each FC's topology, we also compare it with a Clos network generated using roughly the same number of switches with throughput optimized (see Appendix A.3).

We compute throughput values under all-to-all traffic patterns, uniform random traffic patterns and near-worst traffic patterns. In an all-to-all pattern, each server sends an equal amount of traffic to all other servers. In a uniform random pattern, each ToR randomly picks 10% of ToRs to communicate. To generate near-worst patterns, we 1) first construct a complete bipartite graph  $B$  with  $N$  source nodes and  $N$  destination nodes, where the weight of the edge  $(s, d)$  is the length of the shortest path from ToR  $s$  to ToR  $d$ ; 2) and then find the permutation matrix with the maximum weight. This approach was also adopted in [21, 28] to generate near-worst patterns.

**Table 2:** Using Virtual-Layered Cabling Has Little Impact on FC’s Routing Statistics.

Number of Switches	Number of Servers	$k$	Port Count of Virtual Switches	$\eta$	Cabling Constraint	Average Number of Paths	Average Path Length	Minimum Number of Paths
500	12000	4	[7, 13, 13, 7]	4	N	16.08	4.57	12.00
					Y	16.10	4.57	12.00
1000	24000	4	[7, 13, 13, 7]	7	N	14.01	4.86	9.00
					Y	14.01	4.86	9.00
2000	48000	5	[5, 10, 10, 10, 5]	13	N	15.77	5.36	11.00
					Y	15.77	5.36	11.00
3000	72000	5	[5, 10, 10, 10, 5]	19	N	14.66	5.54	10.00
					Y	14.66	5.53	10.00
5000	120000	5	[5, 10, 10, 10, 5]	32	N	13.28	5.75	9.00
					Y	13.28	5.75	9.00



(a) Throughput of the all to all traffic matrix. (The ECMP and KSP curves overlap together.) (b) Throughput of uniform random traffic matrices (averaged over 10 runs). (c) Throughput of the near-worst permutation traffic matrix.

**Figure 7:** Throughput Simulation Results under ECMP, EDST, Edge-disjoint Virtual Up-down and 32-way KSP routing.

#### 4.1 FC’s Routing vs EDST Routing

The EDST routing is CBD-free for expander graphs. A random  $s$ -regular graph has  $s/2$  edge-disjoint spanning trees with high probability [29]. Thus, EDST is a direct competitor of the Edge-disjoint Virtual Up-down Routing for FC’s topology.

As shown in Fig. 10, FC’s edge-disjoint virtual up-down routing (denoted by “DISJOINT UP-DOWN”) performs consistently better than EDST for all the traffic patterns. When the network is small ( $N = 50$ ), FC’s routing achieves  $2 \times$  throughput of the EDST routing. As the network size increases, the performance of the EDST routing deteriorates quickly. When  $N = 500$ , the performance gain of FC’s routing becomes  $10 \times$  and the gain keeps increasing with the network scale.

There are two reasons that lead to the poor performance of the EDST routing. First, existing edge-disjoint spanning tree (EDST) algorithms [34, 35] can find the maximum number of spanning trees, but there is no guarantee that the height of each spanning tree found is small. When we perform routing in a *tall* spanning tree, the average hop count would be large. This is also justified in the following routing-path analysis. Second, some links remain unused in the EDST routing. In an expander graph with  $N$  ToR switches, each spanning tree contains  $N - 1$  links. Note that the total number of ToR-to-ToR links is  $Ns/2$ . When  $Ns/2$  is not divisible by  $N - 1$ , there must be links not used by any spanning tree.

**Routing-Path Analysis:** For several FC’s topologies of different sizes ( $N = 50/100/200/300/500$ ), we analyze the routing paths under FC’s routing and the EDST routing. We calculate three metrics, including average number of paths, average length of paths and average length of the shortest paths. As

shown in Table 3, although the EDST routing could find more paths than FC’s routing, its average path length is much higher. When  $N = 50$ , the average path length under FC’s routing is  $1 - 3.86/7.69 \approx 50\%$  lower than that under the EDST routing. As  $N$  increases to 500, the reduction of average path length becomes  $1 - 5.22/21.96 \approx 76\%$ . We expect that this number will continue to increase for larger networks. The EDST routing cannot guarantee a small routing path length. In contrast, the parameter  $k$  restricts that FC’s routing path length cannot exceed  $2k$  and  $k$  increases slowly with  $N$ .

#### 4.2 FC’s Routing vs ECMP/KSP Routing

ECMP/KSP are widely-used routing protocols for expander graphs. In FC’s topology, ECMP’s throughput fluctuates significantly because ECMP cannot provide enough path diversity; KSP’s throughput is more stable under different traffic patterns. This coincides with the findings in Jellyfish [38].

We note that KSP’s throughput is consistently higher than that of the FC’s edge-disjoint virtual up-down routing. However, deploying KSP routing in expander networks poses a deadlock risk. We have shown in Section 2.2.1.1 that the probability that ECMP/KSP routing contains CBDs is close to 1. Although containing CBDs is not sufficient to trigger deadlocks, we will show in Section 5.1 that ECMP/KSP could indeed trigger deadlocks in certain cases in a real testbed.

#### 4.3 FC vs Clos

Clos is the de facto standard topology for data centers and has witnessed the successful deployment of RDMA in produc-

**Table 3:** Edge-Disjoint Virtual Up-down Routing vs. the EDST Routing.

Number of Switches	Number of servers	k	Port Count of Virtual Switches	Routing	Average Number of Paths	Average Path Length	Average Shortest Path Length
50	700	4	[3, 6, 6, 3]	Edge Disjoint Up-down	8.02	3.86	2.68
				EDST	9.00	7.69	3.12
100	1400	4	[3, 6, 6, 3]	Edge Disjoint Up-down	6.43	4.22	3.04
				EDST	9.00	10.01	4.04
200	2800	4	[3, 6, 6, 3]	Edge Disjoint Up-down	4.99	4.56	3.44
				EDST	9.00	14.01	5.50
300	4200	4	[3, 6, 6, 3]	Edge Disjoint Up-down	4.24	4.75	3.70
				EDST	9.00	16.70	6.52
500	7000	5	[2, 4, 4, 5, 3]	Edge Disjoint Up-down	4.55	5.22	4.00
				EDST	9.00	21.96	8.14

tion [15]. To ensure fair comparisons, given an FC’s topology with  $N$  ToR switches and  $H$  servers, we choose a Clos network that offers the maximum throughput to the  $H$  servers using roughly the same number of switches.

As shown in Fig. 13(a) and 13(b), FC attains  $1.1 - 2 \times$  the throughput of Clos networks. Note that there is a decrease in throughput when the network size changes from 700 to 1400. The reason is that constructing a two-layered Clos becomes infeasible as the number of servers increases to 1400.

However, under near-worst traffic patterns, FC’s throughput can be 15% – 50% lower than that of the Clos networks. We argue that this issue can be resolved when a layer of OCSs is used to interconnect different ToRs. If the real traffic pattern is close to a near-worst pattern of the current topology, we can reconfigure the OCSs to generate a new FC’s topology.

**Remark:** We believe that the above throughput analysis results can truthfully reflect the network performance for RoCEv2 traffic under normal conditions. To validate this statement, We perform a packet-level simulation for a medium-scale data center with 1152 hosts in Appendix A.4.3, and find the results match our throughput analysis.

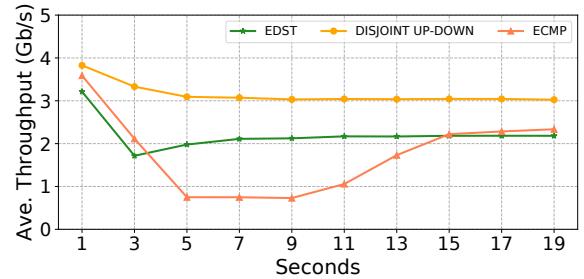
## 5 Formation and Impact of Deadlocks

We study how to trigger deadlocks and understand the impact of deadlocks via both testbed experiments and simulations. We will show that under extreme but practical cases, FC’s edge-disjoint virtual up-down routing is still deadlock-free; but ECMP/KSP could trigger deadlocks.

### 5.1 Trigger Deadlocks in a Real Testbed

We build a small testbed using four switches, each with 8 50Gbps ports. (The four switches are virtualized from a single CE12800 switch. The original port speed is 100Gbps and we limit the port speed as 50Gbps.) This testbed has 16 servers, each equipped with one Mellanox CX5 NIC with maximum rate configured as 50Gbps. (We use PCIE-3.0  $\times 8$  to connect to the NICs, and thus these NICs cannot run at a rate higher than 64Gbps.) Each switch in this testbed has four ports connected to other switches and four ports connected to four servers.

We virtually split each switch into 3 virtual switches, with 1, 2, 1 number of ports respectively. The connections between FC’s virtual switches are shown in Fig. 4(b), and the resulting topology is shown in Fig. 4(a). This topology can be also viewed as a subgraph of a large expander graph (see switches A, B, C, D in Fig. 2). If a deadlock occurs in this subgraph, a PFC storm will quickly propagate to the entire network.



**Figure 8:** Average throughput of the testbed experiment.

We implement ECMP, edge-disjoint virtual up-down and EDST routings using ACL rules in our testbed. We enable PFC to guarantee that the network is loss-free. The PFC-pause threshold  $XOFF$  is set to 50KB and the PFC-resume threshold  $XON$  is set to 47KB. Note that these PFC thresholds are lower than the recommended values. This allows the network to trigger more PFC pauses. As we will see shortly, the virtual up-down routing is deadlock-free even in this extreme situation. Note that this setup can be viewed as a misconfiguration of network switches. Microsoft reports that switch misconfiguration accounts for 38% of the high-impact failures in their data centers [43]. In a PFC-storm incident reported also by Microsoft [15], a switch parameter was misconfigured such that PFC PAUSE frames could be triggered more easily.

We generate RoCEv2 traffic using the “ib\_write\_bw [32]” command. For every NIC, we establish an RDMA connection with every NIC under a different ToR switch. For example, NIC1 under the first ToR switch sends traffic to NIC5, NIC6, ..., NIC16. In total, we establish  $16 \times 12 = 192$  RDMA connections. We configure the “-run\_infinity” parameter at the client side of each connection to run the test indefinitely

until interrupted by external.

**Results:** In the first experiment, we apply ECMP routing. ECMP is not CBD-free in this testbed. We see a deadlock after running our testbed for just a few seconds. (KSP typically generates more paths than ECMP, and thus KSP could also trigger deadlocks.) When deadlock happens, a large number of RDMA connections are broken. We deep dive into the source code of “ib\_write\_bw” to understand why many connections are torn down abnormally. We found that the PFC-deadlocks cause the verbs API “ibv\_post\_send” to fail and return an error code to the main program of “ib\_write\_bw”. Once the main program catches the exception code, “ib\_write\_bw” will stop sending traffic and clean up the resources. Note that, if we use dynamic PFC thresholds or use the recommended values to set static PFC thresholds ( $XOFF = 800KB, XON = 797KB$ ), we could not observe PFC deadlocks under ECMP routing. However, this does not eliminate the deadlock risk for ECMP.

In the second and third experiments, we set up the edge-disjoint virtual up-down and the EDST routing respectively to run the same test. In this case, we do not see any deadlock even under low PFC pause/resume thresholds and all RDMA connections can work continuously. This experiment demonstrates that both the virtual up-down routing and the EDST routing can avoid PFC-deadlock in lossless Ethernet.

Finally, we track the average throughput for all the 192 RDMA connections under different routing strategies over one minute, and plot the results in Fig. 8. The virtual up-down routing attains the highest average throughput, which is about 50% higher than that of the EDST routing. Under ECMP routing, the average throughput drops quickly at the first 5 seconds due to PFC deadlocks. Although the average throughput of ECMP increases after deadlock recovery, 66% of the RDMA connections have already failed.

## 5.2 Understanding Deadlocks via Simulation

We perform packet-level simulation using an open-source NS3 project for HPCC [17]. The objective is to understand how a deadlock is triggered.

We use the same testbed topology (Fig. 4(a)) in our simulation. We generate 192 flows at time 0, and set all the flow sizes as 100MB. For different flows, we either disable congestion control or enable DCQCN for congestion control. When DCQCN is enabled, we set the ECN-marking related parameters as  $Kmin = 5KB, Kmax = 200KB, P_{max} = 0.01$  as suggested by the DCQCN paper [48]. To simulate the extreme cases where lots of PFC pauses are triggered, we set a small PFC-pause threshold and a small PFC-resume threshold ( $XOFF = 50KB, XON = 47KB$ ).

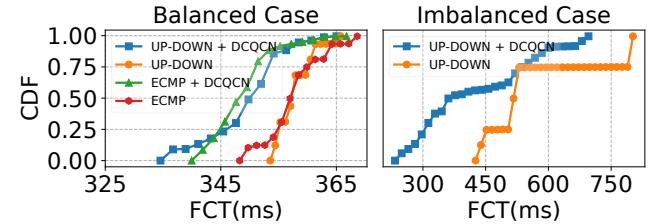
We evaluate ECMP and FC’s edge-disjoint virtual up-down routing. For adjacent switch pairs, there are two paths under both ECMP and FC’s edge-disjoint virtual up-down routing. For non-adjacent switch pairs, there are 8 shortest paths (4 clock-wise paths and 4 counter-clock-wise paths) under

ECMP routing and 2 edge-disjoint virtual up-down paths (1 clock-wise path and 1 counter-clock-wise path) under FC’s routing. We assign a path to each flow using two strategies:

**Balanced Allocation:** There are 16 flows generated between every switch pair and we assign the same number of flows to each path under both routing strategies. In this case, every link between adjacent switch pair is shared by exactly 16 flows.

**Imbalanced Allocation:** Flows between adjacent switch pairs are still equally assigned to all the paths; but flows between non-adjacent switch pairs are only assigned to the clock-wise paths. This situation could happen due to hashing imbalance. In this case, every clock-wise link is shared by 24 flows, which becomes the bottleneck of the network. Incast can thus happen at the 4 switches. In addition, under ECMP routing, the following paths  $\{[e_1, e_2], [e_2, e_3], [e_3, e_4], [e_4, e_1]\}$  form a CBD (actually there are more CBDs), which makes ECMP prone to deadlocks.

**Results:** Under balanced allocation, we do not see deadlocks even if we use a small static PFC threshold and disable DCQCN. In Fig. 9, we compare the CDFs of the FCTs (Flow Completion Time) under both ECMP and FC’s edge-disjoint virtual up-down routing with and without DCQCN. Both routing strategies yield similar FCT performance.



**Figure 9:** CDF of the FCTs of ECMP and FC’s Routing under Balanced/Imbalanced Allocation in the Testbed Topology.

Under imbalanced allocation, FC’s routing can still finish all the flows and the FCT performance is shown in Fig. 9. In contrast, the ECMP routing triggers a deadlock even if we enable DCQCN. To rootcause this issue, we record all the PFC pauses and PFC resumes. We find 4 critical PAUSE signals that lead to the deadlock: 1) at time 531 us,  $S_1$  sends a PAUSE to the link  $e_4$ ; 2) at time 537 us,  $S_4$  sends a PAUSE to the link  $e_3$ ; 3) at time 543 us,  $S_3$  sends a PAUSE to the link  $e_2$ ; 4) at time 552 us,  $S_2$  sends a PAUSE to the link  $e_1$ . These events happen within just 21 us.

**Takeaway:** A DCN suffers from a high risk of deadlocks, when the following three conditions are met: 1) there exist CBDs in the network; 2) links in the CBDs are congested; 3) PFCs are triggered more frequently than usual. If we apply ECMP/KSP routing in an expander graph, we may have to constantly monitor the congested links and the abnormal switch behaviors. FC’s design completely eliminates CBDs, and thus could significantly simplify the RoCEv2 deployment.

## 6 Discussion

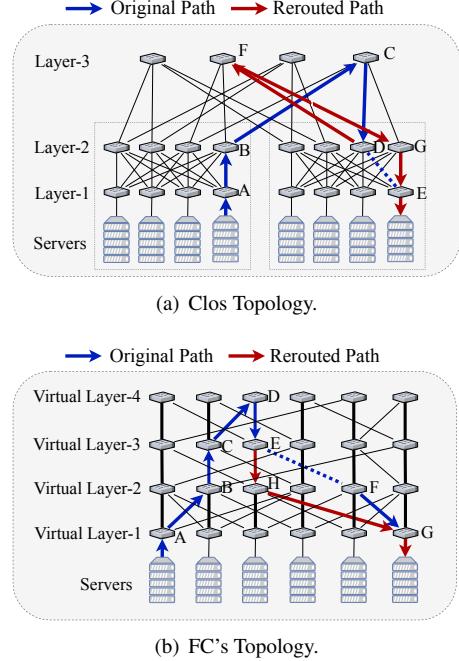
### 6.1 Handling Link/Node Failures

Link/node failures are common in practical data centers [12]. When a link/node fails, to avoid packet drop, local rerouting is performed to forward the affected packets along a different path to the destination [25]. Unfortunately, local rerouting may introduce CBDs and cause deadlocks even if the original network is CBD-free [19]. Consider the Clos network in Fig. 10(a). Initially, packets from ToR A to ToR E follow an up-down path  $A \rightarrow B \rightarrow C \rightarrow D \rightarrow E$ . When the link  $DE$  fails, packets that arrived at the switch  $D$  cannot find an alternative downstream path to  $E$  and thus are bounced back to  $F$ . Then, the path from  $A$  to  $E$  becomes  $A \rightarrow B \rightarrow C \rightarrow D \rightarrow F \rightarrow G \rightarrow E$ . This path contains a down-up bounce, which could potentially introduce CBDs into Clos networks.

To avoid deadlocks in Clos networks under link/node failures, Tagger [19] is proposed to add a tag to all the packets, increasing the tag on the bounce and putting packets with different tags into different lossless queues. This approach should also work for FC because we can treat FC as a virtual multi-layered network. Nevertheless, better approach may exist for FC. In Clos networks, every top-layer switch has a unique path to every ToR switch and thus every packet affected by a downstream link/node failure has to be bounced back to another top-layer switch. In contrast, every packet affected by a link/node failure in FC can freely choose any virtual layer as long as there is a link to forward this packet, because every virtual switch in the same column (see Fig. 10(b)) belongs to the same physical switch. For example, there is a flow from  $A$  to  $G$  in Fig. 10(b) and the original path is  $A \rightarrow B \rightarrow C \rightarrow D \rightarrow E \rightarrow F \rightarrow G$ . When the link  $EF$  fails, the affected packets can be rerouted to  $A \rightarrow B \rightarrow C \rightarrow D \rightarrow E \rightarrow H \rightarrow G$ . This new path is still an up-down path and thus tagging is not required. (Admittedly, if the rerouted path contains a down-up bounce, we still need to update the packet tags.) Based on the above analysis, we suspect that FC could be more efficient in handling link/node failures than a Clos network. We will explore this further in our future work.

### 6.2 Handling Route Reconfiguration

Route reconfiguration is common in data centers, which could happen when 1) new flows join/leave the network; 2) DCN topology changes; 3) Traffic Engineering is enabled; 4) an SDN controller reoptimizes routing paths after link/node failures. FC's design makes it easy to handle route reconfiguration. FC performs virtual up-down routing. As long as the virtual layers remain unchanged (i.e., which ToR port belongs to which virtual layer), the combined set of the original paths and the post-reconfiguration paths is CBD-free. This could



**Figure 10:** Rerouting Under Link Failures.

dramatically simplify the workflow of route reconfiguration, because any transient state during route reconfiguration is guaranteed to be deadlock-free.

In rare cases, e.g., after data center expansion, we may need to change the virtual layers because the original number of layers may not be able to support a larger-scale network. In this case, there could be a CBD in a transient state during route reconfiguration. Existing solutions on deadlock-free route reconfiguration [8, 19, 26, 30] can be applied here.

## 7 Conclusion

We present FC, a topology-routing co-designed methodology to eliminate PFC-induced deadlocks, for cost-effective and safe deployment of RoCEv2 over expander networks. Motivated by the fact that the up-down routing paths of multi-layered Clos networks are CBD-free, we design FC's topology to exhibit a virtual layered structure, and propose an edge-disjoint virtual up-down routing for FC that is guaranteed to be CBD-free. We evaluate FC against several competitors using throughput analysis, testbed implementation and packet-level simulation. Our evaluation results demonstrate that 1) FC is deadlock-free while ECMP/KSP may trigger deadlocks; 2) FC significantly reduces average hop count and improves network throughput over the state-of-art EDST-based routing strategy; 3) FC attains higher throughput than Clos networks built using the same number of switches under all-to-all and uniform random patterns. These properties make FC a promising design for deadlock prevention in expander graphs.

## References

- [1] Mincostflow solver of ortools. [https://developers.google.com/optimization/reference/graph/min\\_cost\\_flow](https://developers.google.com/optimization/reference/graph/min_cost_flow).
- [2] M. Al-Fares, A. Loukissas, and A. Vahdat. A scalable, commodity data center network architecture. In *SIGCOMM*, pages 63–74, 2008.
- [3] H. Ballani, P. Costa, R. Behrendt, D. Cletheroe, I. Haller, K. Jozwik, F. Karinou, S. Lange, K. Shi, B. Thomsen, and H. Williams. Sirius: A flat datacenter network with nanosecond optical switching. In *SIGCOMM*, 2020.
- [4] M. Besta and T. Hoefer. Slim fly: A cost effective low-diameter network topology. In *SC*, pages 348–359. IEEE, 2014.
- [5] M. Besta, M. Schneider, M. Konieczny, K. Cynk, E. Henriksson, S. D. Girolamo, A. Singla, and T. Hoefer. Fat-paths: Routing in supercomputers and data centers when shortest paths fall short. In *SC*, pages 1–18. IEEE, 2020.
- [6] B. Bollobás. The isoperimetric number of random regular graphs. *European Journal of combinatorics*, 9(3):241–244, 1988.
- [7] P. Cao, S. Zhao, M. Y. The, Y. Liu, and X. Wang. Trod: Evolving from electrical data center to optical data center. In *ICNP*, pages 1–11. IEEE, 2021.
- [8] J.-J. Crespo, J. L. Sánchez, F. J. Alfaro-Cortés, J. Flich, and J. Duato. Upr: deadlock-free dynamic network reconfiguration by exploiting channel dependency graph compatibility. *The Journal of Supercomputing*, 77:12826–12856, 2021.
- [9] W. J. Dally and C. L. Seitz. Deadlock-free message routing in multiprocessor interconnection networks. 1988.
- [10] Y. Gao, Q. Li, L. Tang, Y. Xi, P. Zhang, W. Peng, B. Li, Y. Wu, S. Liu, L. Yan, et al. When cloud storage meets rdma. In *NSDI*, pages 519–533, 2021.
- [11] M. Gerla and L. Kleinrock. Flow control: A comparative survey. *IEEE Transactions on Communications*, 28(4):553–574, 1980.
- [12] P. Gill, N. Jain, and N. Nagappan. Understanding network failures in data centers: Measurement, analysis, and implications. In *SIGCOMM*, pages 350–361, 2011.
- [13] A. V. Goldberg and M. Kharitonov. On implementing scaling push-relabel algorithms. *Network Flows and Matching: First DIMACS Implementation Challenge*, 1993.
- [14] A. Greenberg, J. R. Hamilton, N. Jain, S. Kandula, C. Kim, P. Lahiri, D. A. Maltz, P. Patel, and S. Sen-gupta. Vi2: A scalable and flexible data center network. In *SIGCOMM*, pages 51–62, 2009.
- [15] C. Guo, H. Wu, Z. Deng, G. Soni, J. Ye, J. Padhye, and M. Lipshteyn. Rdma over commodity ethernet at scale. In *Proceedings of the 2016 ACM SIGCOMM Conference*, pages 202–215, 2016.
- [16] M. Handley, C. Raiciu, A. Agache, A. Voinescu, A. W. Moore, G. Antichi, and M. Wójcik. Re-architecting datacenter networks and stacks for low latency and high performance. In *SIGCOMM*, pages 29–42, 2017.
- [17] HPCC. <https://github.com/alibaba-edu/High-Precision-Congestion-Control>.
- [18] S. Hu, Y. Zhu, P. Cheng, C. Guo, K. Tan, J. Padhye, and K. Chen. Deadlocks in datacenter networks: Why do they form, and how to avoid them. In *Proceedings of the 15th ACM Workshop on Hot Topics in Networks*, pages 92–98, 2016.
- [19] S. Hu, Y. Zhu, P. Cheng, C. Guo, K. Tan, J. Padhye, and K. Chen. Tagger: Practical pfc deadlock prevention in data center networks. In *Proceedings of the 13th International Conference on emerging Networking EXperiments and Technologies*, pages 451–463, 2017.
- [20] IEEE. <https://1.ieee802.org/dcb/802-1qbb/>.
- [21] S. A. Jyothi, A. Singla, P. B. Godfrey, and A. Kolla. Measuring and understanding throughput of network topologies. In *SC*, 2016.
- [22] M. Karol, S. J. Golestani, and D. Lee. Prevention of deadlocks and livelocks in lossless backpressured packet networks. *IEEE/ACM Transactions on Networking*, 11(6):923–934, 2003.
- [23] S. Kassing, A. Valadarsky, G. Shahaf, M. Schapira, and A. Singla. Beyond fat-trees without antennae, mirrors, and disco-balls. In *SIGCOMM*, pages 281–294, 2017.
- [24] Y. Li, R. Miao, H. H. Liu, Y. Zhuang, F. Feng, L. Tang, Z. Cao, M. Zhang, F. Kelly, M. Alizadeh, et al. Hpcc: high precision congestion control. In *SIGCOMM*, pages 44–58, 2019.
- [25] V. Liu, D. Halperin, A. Krishnamurthy, and T. Anderson. F10: A fault-tolerant engineered network. In *NSDI*, 2013.
- [26] O. Lysne, T. M. Pinkston, and J. Duato. A methodology for developing deadlock-free dynamic network reconfiguration processes. part ii. *IEEE Transactions on Parallel and Distributed Systems*, 16(5):428–443, 2005.

- [27] R. Mittal, A. Shpiner, A. Panda, E. Zahavi, A. Krishnamurthy, S. Ratnasamy, and S. Shenker. Revisiting network support for rdma. In *SIGCOMM*, pages 313–326, 2018.
- [28] P. Namyar, S. Supittayapornpong, M. Zhang, M. Yu, and R. Govindan. A throughput-centric view of the performance of datacenter topologies. In *SIGCOMM*, pages 349–369, 2021.
- [29] E. Palmer. On the spanning tree packing number of a graph: A survey. *Discrete Mathematics*, 2001.
- [30] T. M. Pinkston, R. Pang, and J. Duato. Deadlock-free dynamic reconfiguration schemes for increased network dependability. *IEEE Transactions on Parallel and Distributed Systems*, 14(8):780–794, 2003.
- [31] L. Poutievski, O. Mashayekhi, J. Ong, A. Singh, M. Tariq, R. Wang, J. Zhang, V. Beauregard, P. Conner, S. Gribble, et al. Jupiter evolving: transforming google’s datacenter network via optical circuit switches and software-defined networking. In *SIGCOMM*, pages 66–85, 2022.
- [32] L. Rdma. <https://github.com/linux-rdma/perf-test>.
- [33] R. Recio, B. Metzler, P. Culley, J. Hilland, and D. Garcia. A remote direct memory access protocol specification. Technical report, RFC 5040, October, 2007.
- [34] J. Roskind. *Application of Edge Disjoint Trees to Failure Recovery in Data Communication Networks*. PhD thesis, PhD thesis, Department of Electrical Engineering and Computer Science, 1983.
- [35] J. Roskind and R. E. Tarjan. A note on finding minimum-cost edge-disjoint spanning trees. *Mathematics of Operations Research*, 10(4):701–708, 1985.
- [36] A. Roy, H. Zeng, J. Bagga, G. Porter, and A. C. Snoeren. Inside the social network’s (datacenter) network. In *SIGCOMM*, 2015.
- [37] A. Singh, J. Ong, A. Agarwal, G. Anderson, A. Armistead, R. Bannon, S. Boving, G. Desai, B. Felderman, P. Germano, et al. Jupiter rising: A decade of clos topologies and centralized control in google’s datacenter network. pages 183–197, 2015.
- [38] A. Singla, C.-Y. Hong, L. Popa, and P. B. Godfrey. Jellyfish: Networking data centers randomly. In *NSDI*, 2012.
- [39] B. Stephens and A. L. Cox. Deadlock-free local fast failover for arbitrary data center networks. In *INFOCOM*, pages 1–9. IEEE, 2016.
- [40] B. Stephens, A. L. Cox, A. Singla, J. Carter, C. Dixon, and W. Felter. Practical dcb for improved data center networks. In *INFOCOM*, pages 1824–1832. IEEE, 2014.
- [41] C. Technologies. <https://www.calient.net/products/s-series-photonic-switch/>.
- [42] A. Valadarsky, G. Shahaf, M. Dinitz, and M. Schapira. Xpander: Towards optimal-performance datacenters. In *Proceedings of the 12th International Conference on Emerging Networking EXperiments and Technologies*, pages 205–219, 2016.
- [43] X. Wu, D. Turner, C.-C. Chen, D. A. Maltz, X. Yang, L. Yuan, and M. Zhang. Netpilot: Automating datacenter network failure mitigation. In *SIGCOMM*, pages 419–430, 2012.
- [44] J. Y. Yen. Finding the k shortest loopless paths in a network. *management Science*, 17(11):712–716, 1971.
- [45] M. Zhang, R. N. Mysore, S. Supittayapornpong, and R. Govindan. Understanding lifecycle management complexity of datacenter topologies. In *NSDI*, 2019.
- [46] S. Zhao, P. Cao, and X. Wang. Understanding the performance guarantee of physical topology design for optical circuit switched data centers. In *SIGMETRICS*, 2021.
- [47] S. Zhao, R. Wang, J. Zhou, J. Ong, J. C. Mogul, and A. Vahdat. Minimal rewiring: Efficient live expansion for clos data center networks. In *NSDI*, 2019.
- [48] Y. Zhu, Y. Zhu, H. Eran, D. Firestone, D. Firestone, C. Guo, M. Lipshteyn, Y. Liron, J. Padhye, S. Raindel, M. H. Yahia, M. Zhang, and J. Padhye. Congestion control for large-scale rdma deployments. In *SIGCOMM*, 2015.

## A Appendix

### A.1 Finding Edge-Disjoint Paths Using Min-Cost Max-Flow

**Definition 2** (*Min-Cost Max-Flow Problem*) Given a flow network  $G(V, E)$  with

- $u(v, w)$ , upper bound on flow from node  $v$  to node  $w$ ;
- $c(v, w)$ , cost of a unit of flow on  $(v, w)$ ,

and a source-destination pair  $(s, t)$ ,  $[f(v, w)]_{(v,w) \in E}$  is called a flow assignment from  $s$  to  $t$  if the following constraints are met:

1. Capacity constraints:  $0 \leq f(v, w) \leq u(v, w)$ ;
2. Flow conservation constraints:  $\sum_u f(u, v) = \sum_w f(v, w)$  for any node  $v \neq s, t$  and  $\sum_w f(s, w) = \sum_u f(u, t) = F$ . Here  $F$  is called the total amount of flow from  $s$  to  $t$ .

The objective of the min-cost max-flow problem is to find a flow assignment  $[f(v, w)]_{(v,w) \in E}$  with the maximum flow that minimizes

$$\sum_{(v,w)} c(v, w) \cdot f(v, w).$$

Note that the constant parameters  $u(v, w)$  are all positive and  $c(v, w)$  can be either positive or negative. In addition, the min-cost max-flow problem has a very nice property that guarantees integer solutions:

**Theorem 6** (*Integral Flow Theorem*) Given a min-cost max-flow problem, if  $u(v, w)$ 's are all integers, then there exists an integer solution, i.e.,  $f(v, w)$ 's are all integers, such that  $[f(v, w)]_{(v,w) \in E}$  attains the maximum flow with minimum cost.

In fact, when we solve a min-cost max-flow problem with integer bounds using the Scaling Push-Relabel algorithm [1, 13], the resulting optimal solution is guaranteed to be an integer solution.

**Finding Edge-Disjoint Paths:** As a consequence of Theorem 6, we can find the maximum number of edge-disjoint paths from  $s$  to  $t$  using min-cost max-flow. Specifically, let  $E_0$  be the set of links that can be used at most once (see the solid links in Fig. 4(c)), and  $E \setminus E_0$  be the set of links that can be used multiple times (see the dashed links in Fig. 4(c)). If we set the upper bound as  $u(v, w) = 1$  for all the links  $(v, w) \in E_0$  and set the upper bound as  $u(v, w) = \infty$  for all the links  $(v, w) \in E \setminus E_0$ , then the resulting min-cost max-flow solution  $[f(v, w)]_{(v,w) \in E}$  can be decomposed into  $F$  ( $F$  is the maximum flow) paths where links in  $E_0$  can be used at most once. The  $F$  paths can be found by performing Depth First Search  $F$  times (see Algorithm 1). Note that when we perform DFS in line 5 of Algorithm 1, we will never encounter a cycle. Otherwise, by removing this cycle we could obtain another flow assignment with lower cost. Having this observation could slightly simplify the DFS implementation. We do not need to track the set of visited nodes during the DFS search.

---

**Algorithm 1:** Find Edge-Disjoint Paths in the Directed Virtual Up-Down Graph

---

**Input :** A directed virtual up-down graph (see Fig. 4(c)) and a source-destination pair  $(s, t)$ .  
**Output :** Maximum number of edge-disjoint up-down paths from  $s$  to  $t$ .

- 1 Let  $E_0$  be the set of solid lines in the directed virtual up-down graph. Construct a flow graph by setting the link capacity and the link cost as 1 for all links in  $E_0$ , and setting the link capacity as  $\infty$  and the link cost as  $\epsilon$  (an infinitesimal value) for all links not in  $E_0$ .
- 2 Solve the min-cost max-flow problem. Let  $[f(v, w)]_{(v,w) \in E}$  be the optimal solution and let  $F$  be the maximum flow from  $s$  to  $t$ .
- 3 Use  $\mathcal{P}$  to store the set of paths, and initialize  $\mathcal{P} = \emptyset$ .
- 4 **for**  $i$  **in**  $\{1, 2, \dots, F\}$  **do**
- 5     Use Depth First Search to find a path  $P$  from  $s$  to  $t$  such that  $f(e) \geq 1$  for every edge  $e$  in  $P$ .
- 6     Store  $P$  in  $\mathcal{P}$ .
- 7     For every edge  $e$  in  $P$ , decrement  $f(e)$  by one.
- 8 **end**
- 9 Return  $\mathcal{P}$ .

---

### A.2 A Sufficient and Necessary Condition for CBD-Free Routing

We first introduce the concept of *link dependency graph*.

**Definition 3** Given a network  $G(V, E)$  and a path set  $\mathcal{P} = \{P_1, P_2, \dots, P_K\}$ , a link dependency graph  $G'(V', E')$  can be constructed as follows:

1.  $V'$  is the set of directed links used by at least one path  $P \in \mathcal{P}$ ;
2. For any  $e_1, e_2 \in V'$ , there is a directed link from  $e_1$  to  $e_2$  in  $E'$  if and only if  $e_1$  is the next hop of  $e_2$  in one path  $P \in \mathcal{P}$ .

Then, the following theorem offers a sufficient and necessary condition for a set of paths to be CBD-free.

**Theorem 7** Given a network  $G(V, E)$ , a path set  $\mathcal{P} = \{P_1, P_2, \dots, P_K\}$  is CBD free if and only if the corresponding link dependency graph  $G'(V', E')$  contains no loops.

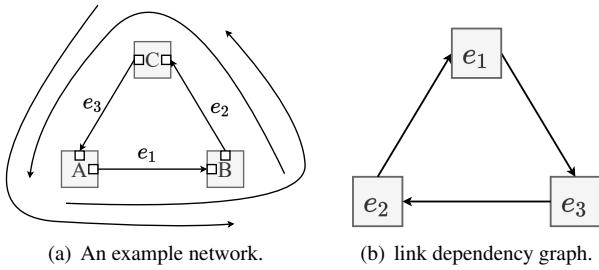
**Proof 2 Necessity  $\Rightarrow$ :** If the path set  $\mathcal{P}$  is CBD free, we prove that  $G'(V', E')$  contains no loops. We prove this by contradiction. Suppose that  $G'(V', E')$  contains a loop  $v'_1 \rightarrow v'_2 \rightarrow \dots \rightarrow v'_s \rightarrow v'_1$ . Let  $e_i$  be the link in  $G(V, E)$  that corresponds to  $v'_i$ . Since  $e_1$  is the next hop of  $e_2$  in a path, if  $e_1$  is paused,  $e_2$  will be paused. Based on the same argument,  $e_3, \dots, e_s$  will be paused. Since  $e_s$  is the next hop of  $e_1$  in a path, the pause of  $e_s$  will in turn pause  $e_1$ . Then, a CBD is formed, which contradicts the assumption that the path set  $\mathcal{P}$  is CBD-free.

**Sufficiency**  $\Leftarrow$ : If  $G'(V', E')$  contains no loops, we prove that the path set  $\mathcal{P}$  is CBD free. We again prove this by contradiction. Suppose that  $\mathcal{P}$  contains a CBD. Then, there must exist a sequence of links  $e_1, e_2, \dots, e_s$  such that  $e_i$  is the next hop of  $e_{i+1}$  in a path and  $e_s$  is the next hop of  $e_1$  in a path. Then, the corresponding vertices of  $e_1, e_2, \dots, e_s$  in  $G'(V', E')$  forms a loop, which contradicts to the assumption that  $G'(V', E')$  contains no loop.

According to Theorem 7, we design Algorithm 2 to check if a set of paths is deadlock-free.

**Algorithm 2:** Check if a set of paths is deadlock-free

- Input** : A set of paths  $\mathcal{P} = \{P_1, P_2, \dots, P_K\}$  and a network  $G(V, E)$ .  
**Output** : Whether  $\mathcal{P}$  is deadlock-free.  
1 Construct a link dependency graph  $G'(V', E')$  based on Definition 3.  
2 Use deep first search to check if  $G'(V', E')$  has a loop.  
3 Return true if  $G'$  has no loop; return false otherwise.



**Figure 11:** Deadlock detection with a link dependency graph.

We use the example in Fig. 11 to illustrate the idea of the deadlock detection algorithm. Given a path set  $\mathcal{P} = \{A \rightarrow B \rightarrow C, B \rightarrow C \rightarrow A, C \rightarrow A \rightarrow B\}$ , we can construct a link dependency graph with three vertices:  $e_1(A \rightarrow B)$ ,  $e_2(B \rightarrow C)$ ,  $e_3(C \rightarrow A)$ . It is easy to see that this link dependency graph contains a loop. Thus, the path set  $\mathcal{P}$  has the risk of deadlock.

Using Theorem 7, we can prove that up-down routing is CBD-free in a multi-layered network.

**Theorem 8** In a multi-layered network, the path set generated by up-down routing is CBD-free.

**Proof 3** For all the links in a multi-layered network, we can define a partial order as follows. A link  $e_1$  is considered smaller than another link  $e_2$  if either of the following three conditions is met:

1.  $e_1$  is an up link while  $e_2$  is a down link;
2. both  $e_1$  and  $e_2$  are down links and  $e_1$  is at a higher layer than  $e_2$ ;

3. both  $e_1$  and  $e_2$  are up links and  $e_1$  is at a lower layer than  $e_2$ .

Then, when we construct a link dependency graph based on up-down paths, it is easy to verify the following fact: if there is a directed link from  $e_1$  to  $e_2$ , we must have  $e_2 < e_1$ . Therefore, the link dependency graph cannot contain a loop. As a result, the path set generated by up-down routing is CBD-free.

### A.2.1 Proof of Lemma 4

**Proof 4** Since the path set  $\mathcal{P}$  is CBD free in  $G(V, E)$ , the corresponding link dependency graph  $G'(V', E')$  must contain no loop. In this case, we could construct a new graph  $G''(V', E'')$  by adding a link from  $v_1 \in V'$  to  $v_k \in V'$  whenever there exists a sequence of node  $v_2, v_3, \dots, v_{k-1} \in V'$  such that  $(v_i, v_{i+1}) \in E'$  for every  $i = 1, 2, \dots, k-1$ . It is easy to check that  $G''(V', E'')$  is also loop-free.

Now we consider the contraction process. Let  $\hat{G}'(\hat{V}', \hat{E}')$  be the link dependency graph of  $(\hat{G}(\hat{V}, \hat{E}), \hat{\mathcal{P}})$ . We can prove that  $\hat{G}'(\hat{V}', \hat{E}')$  is a subgraph of  $G''(V', E'')$ . First, in  $\hat{G}(\hat{V}, \hat{E})$ , the edges within each vertex set  $V_i$  ( $i=1, 2, \dots, m$ ) are removed. Thus,  $\hat{V}' \subseteq V'$ . Second, for any edge  $(e_1, e_2) \in \hat{E}'$ , there must be a path  $\hat{P} \in \hat{\mathcal{P}}$ , such that  $e_1$  is the next hop of  $e_2$  in  $\hat{P}$ . Note that  $\hat{P}$  is obtained by contracting a path  $P \in \mathcal{P}$ . We must have  $e_1$  as a down-streaming hop (not necessarily next hop) of  $e_2$  in  $P$ . Based on the construction of  $G''(V', E'')$ , we know that  $(e_1, e_2) \in V''$ . Therefore,  $\hat{E}' \subseteq V''$ . Based on the above analysis, we immediately know that  $\hat{G}'(\hat{V}', \hat{E}')$  is a subgraph of  $G''(V', E'')$ . Since  $G''(V', E'')$  is loop-free,  $\hat{G}'(\hat{V}', \hat{E}')$  must also be loop-free. Then, according to Theorem 7, we must have that the path set  $\hat{\mathcal{P}}$  is CBD free in the topology  $\hat{G}(\hat{V}, \hat{E})$ .

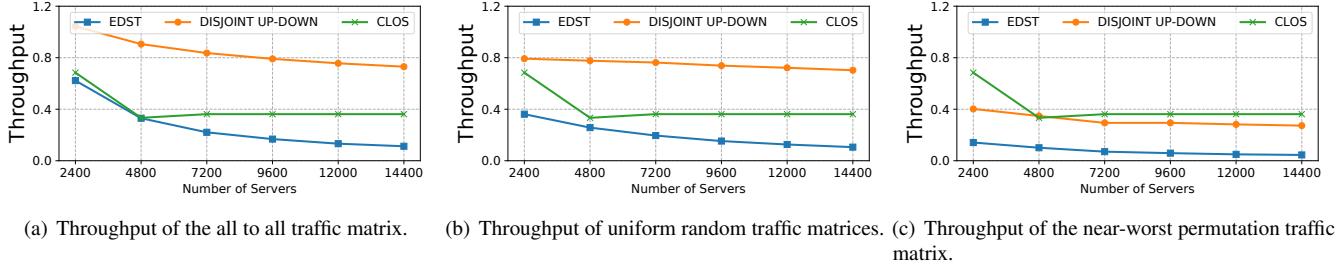
### A.3 Generating a Clos Network with $H$ Hosts Using $N$ $p$ -Port Switches

Given  $N$   $p$ -port switches and  $H$  hosts, we study how to construct a Clos network with the maximum throughput.

We first consider a 2-layered Clos Network. For each switch, let  $h$  be the number of ports connected to hosts. Then, the total number of switches in the first layer (i.e., the ToR layer) is  $\lceil H/h \rceil$ . As long as  $\lceil H/h \rceil \leq p$ , we can put  $p-h$  switches in the second layer and create a complete bipartite graph between the ToR switches and the switches in the second layer. In total,  $\lceil H/h \rceil + p-h$  switches are used. To maximize throughput, we only need to find the smallest  $h$  by solving the following optimization problem:

$$\min h \text{ such that } \lceil H/h \rceil \leq p, \lceil H/h \rceil + p - h \leq N. \quad (5)$$

In many cases, it may not be feasible to construct a 2-layered Clos network or a 2-layered Clos network may not be throughput optimal. Hence, we also need to study how to construct a multi-layered Clos network.



**Figure 12:** Throughput Simulation Results Using 64-port Switches.

We adopt a trial-and-error approach to find the throughput optimal  $L$ -layered Clos network ( $L = 3, 4, \dots$ ). Starting from  $h = 1$ , we try if it is possible to construct an  $L$ -layered Clos network using at most  $N$  switches. If it is possible, we obtain the optimal  $h$  for the  $L$ -layered Clos network; otherwise, we increase  $h$  by one and retry the construction.

Starting from  $L = 2$ , we could use the above approach to find the best  $h(L)$  for every  $L$  ( $h(L) = \infty$  if it is not feasible to construct an  $L$ -layered Clos network).  $h(L)$  may decrease at the beginning, but will eventually increase with respect to  $L$ . Whenever we see  $h(L) < h(L+1)$ , we can stop and return the minimum value of  $h$ , denoted by  $h^*$ . With  $h^*$ , the optimal throughput is  $(p - h^*)/h^*$ . When  $h^* \leq \lfloor p/2 \rfloor$ , the optimal throughput becomes larger than 1. In this case, the DCN offers abundant capacity while the access links between servers and ToRs become the bottleneck.

## A.4 Additional Results

### A.4.1 Throughput Analysis

Clos, FC and Expander+EDST are three network architectures that are guaranteed to be deadlock-free. For networks built using 32-port switches, we demonstrate in Section 4 that

1. FC consistently outperforms Expander+EDST;
2. FC achieves higher throughput than Clos networks under all-to-all and uniform random traffic patterns.

In this section, we generate FC’s topologies of different sizes using up to 600 64-port ToR switches. Each ToR switch has 40 ports connected to other switches and 24 ports connected to servers. The number of virtual layers  $k$  is chosen based on the strategy (\*) in Section 3.2.3. We evaluate both FC’s edge-disjoint virtual up-down routing and the EDST routing. For each FC’s topology, we also compare it with a Clos network generated using roughly the same number of switches with throughput optimized. From Fig. 12, we can see that the above conclusions on FC’s throughput benefits also hold for networks built using 64-port switches.

### A.4.2 Routing-Path Analysis

We then perform the same routing-path analysis for FC’s edge-disjoint virtual up-down routing and the EDST routing. We generate FC’s topologies of different sizes ( $N = 500/1000/2000/3000$ ) using 64-port ToR switches with  $s = 40$ . As shown in Table 4, the average path length under FC’s routing is still much shorter than that under the EDST routing.

### A.4.3 Packet-Level Simulation

We cross-validate our throughput analysis using a packet-level simulator. We generate an FC’s topology using 144 32-port switches. Each switch has 8 ports connected to hosts and 24 ports connected to other switches. In total, there are 1152 hosts. On top of this topology, we run FC’s routing or the EDST routing. We also generate a Clos network using 148 32-port switches with throughput optimized. This Clos network has 64 ToR switches, 56 aggregation switches and 28 spine switches. Each ToR has 18 ports connected to hosts and 14 ports connected to the aggregation switches. The total number of hosts is still 1152. For this Clos topology, we use up-down routing. The port speed is set as 25Gbps.

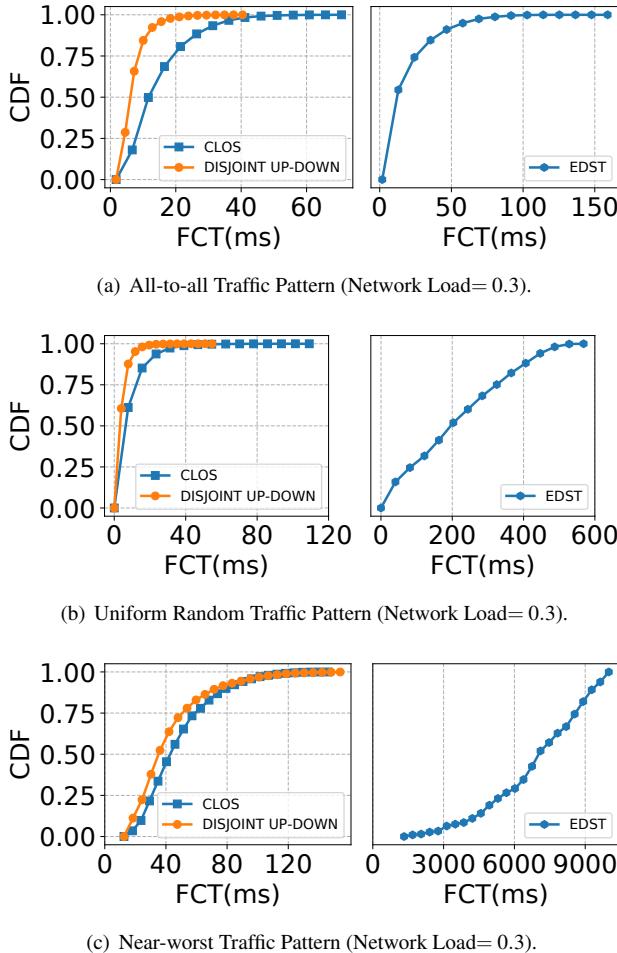
We generate three sets of flows based on the all-to-all traffic pattern, a uniform random traffic pattern (each ToR choose 12.5% of ToRs to communicate) and the near-worst traffic pattern. The default network load is 0.3, meaning that the maximum ingress/egress traffic of each ToR is  $0.3 \times \text{Number of Hosts per ToR} \times 25\text{Gbps}$ . For all the flows, we enable DCQCN for congestion control. We adopt dynamic PFC threshold such that the PFC is triggered when an ingress queue consumes more than 11% of the free switch buffer as suggested by HPCC [24]. We evaluate performance based on the flow completion time (FCT).

We plot the CDFs for FCTs in Fig. 13. Apparently, the FCT performance under FC’s routing is much better than that under the EDST routing. Hence, we mainly focus on the comparison between FC and Clos.

Under the all-to-all traffic pattern and the uniform random traffic pattern, FC achieves clearly better FCT performance than Clos because it has higher throughput. This coincides with our throughput analysis in Section 4.3.

**Table 4:** Edge-Disjoint Virtual Up-down Routing vs. the EDST Routing.

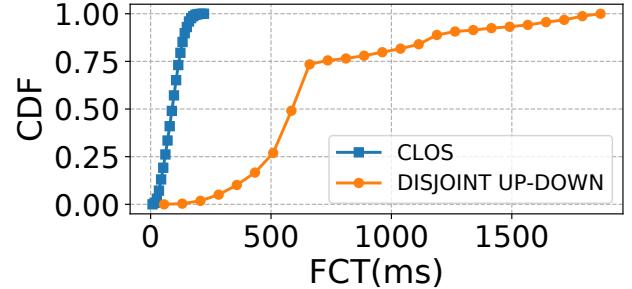
Number of Switches	Number of servers	k	Port Count of Virtual Switches	Routing	Average Number of Paths	Average Path Length	Average Shortest Path Length
500	12000	4	[7, 13, 13, 7]	Edge Disjoint Up-down	16.08	4.57	3.20
				EDST	20	18.34	5.26
1000	24000	4	[7, 13, 13, 7]	Edge Disjoint Up-down	14.01	4.86	3.50
				EDST	20	26.52	6.99
2000	48000	4	[7, 13, 13, 7]	Edge Disjoint Up-down	11.85	5.13	7.00
				EDST	20	33.49	9.12
3000	72000	4	[7, 13, 13, 7]	Edge Disjoint Up-down	10.63	5.30	6.00
				EDST	20	39.82	10.45



**Figure 13:** Compare FCTs for FC, Clos and Exander+EDST.

However, under the near-worst traffic pattern, we find that FC's FCT performance is just slightly worse than the Clos network's FCT performance. In contrast, our throughput analysis in Section 4.3 suggests that FC's throughput is lower than the corresponding Clos network's throughput. (In this case, FC's throughput under the near-worst pattern is about 0.5, while the Clos network's throughput is about 0.78.) The reason is

that, the average hop count under FC is shorter than that under a Clos network; when the network is not congested, having a smaller average hop count compensates for the throughput gap between FC and Clos. Nevertheless, as network load increases, FC will encounter more severe congestion than Clos. We perform another simulation for the near-worst traffic pattern with network load increased from 0.3 to 0.7. (More specifically, we increase the size of each flow by 7/3 times.) As shown in Fig. 14, we can see that FC's FCT performance is much worse than the Clos network's FCT performance. In fact, we see a large amount of PFC PAUSE frames in FC's network. But the good news is, there is no deadlock.



**Figure 14:** FC vs Clos under Near-Worst Traffic Pattern. (Network Load= 0.7)