# MAESTRO: MAke-bEfore-break StraTegy for Reconfiguration in Optical Datacenters

**Sandeep Kumar Singh[1], Che-Yu Liu[2], Roberto Proietti[3], and S. J. Ben Yoo[1]**

*[1]Department of Electrical and Computer Engineering, University of California, Davis, 95616, CA, USA*
*[2]Department of Computer Science, University of California, Davis, 95616, CA, USA*
*[3]Dipartimento di Elettronica e Telecomunicazioni, Politecnico di Torino, 10129, Torino, Italy*

{*sansingh,cyliu,sbyoo*}*@ucdavis.edu; roberto.proietti@polito.it*

**Abstract:** We present a MAke-bEfore-break StraTegy for Reconfiguration in Optical datacenters (MAESTRO). The simulation results show a reduction in packet loss by up to 98% compared to a baseline reconfiguration method. © 2023 The Author(s)

## 1. Introduction

Reconfigurable datacenters (DCs) and high-performing computing (HPC) systems can deliver higher speed, capacity, and energy efficiency when dealing with skewed communication-intensive workloads (e.g., distributed deep learning) with low latency requirements [1, 2]. While academia has been looking into the potential benefits of topology and bandwidth reconfiguration for a while, Google revealed recently that their Jupyter datacenter (DC) incorporates MEMS-based optical circuit switches (OCSs) in the spine layer to enable reconfigurability between aggregation switches [3]. This allows to match the estimated or predicted traffic by steering more bandwidth (or wavelengths) to the hot-spot fiber-links in a dynamic traffic scenario [1, 4]. However, a loss-free or hitless reconfiguration on live network fabric is a challenging task. It might compromise link availability, requiring a topology and traffic engineering strategy to avoid or minimize traffic disruption.

This work presents a MAke-bEfore-break StraTegy for Reconfiguration in Optical datacenters (MAESTRO). We devise a combination of topology and traffic engineering during the reconfiguration phase. MAESTRO routes traffic flows on a residual topology while draining and switching the logical links, i.e., wavelengths. After updating the logical topology, it reroutes the traffic flows to exploit the steered bandwidth. We evaluate the performance of MAESTRO against a baseline optical switch reconfiguration (OSR) method. OSR breaks down the logical links to be reconfigured for the whole reconfiguration duration and reroutes traffic flows over the updated topology. We evaluate both reconfiguration methods on two photonic switches-enabled Hyper-X networks [2]. Our simulation results show up to 98% improvement in packet loss by the MAESTRO scheme compared to the OSR approach.
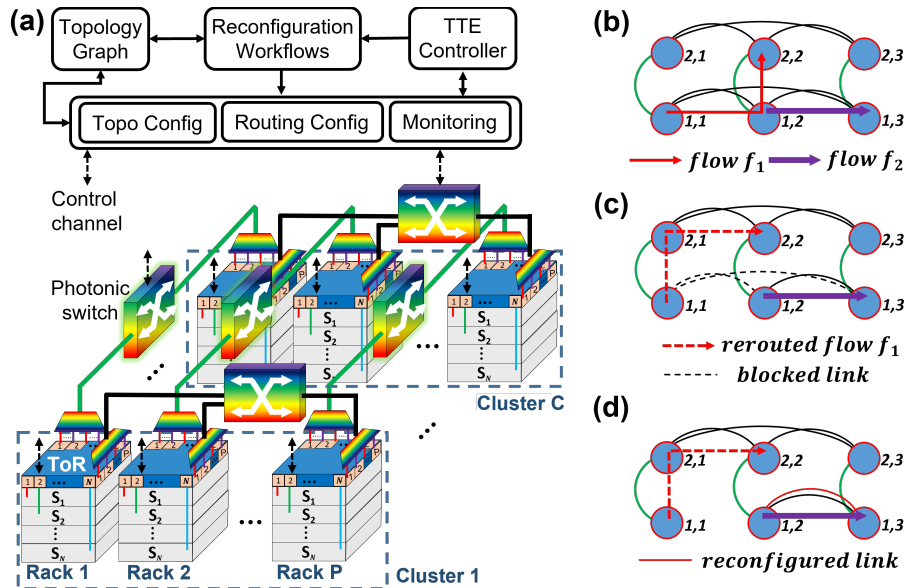


Fig. 1. (a) An architecture of a 2D-Hyper-X DC/HPC system with $N$ servers per rack. ToR switches in each row and column are interconnected with an optical switch. (b-d) MAESTRO operations.

---

**Algorithm 1** MAESTRO procedure

---

1: **Given**: Demand matrix $\widehat{\mathbf{D}}$, current topology graph $G$, connectivity graph $G_0$.
2: Compute matching topology $G^+$ for the demand matrix $\widehat{\mathbf{D}}$ considering the physical connectivity in $G_0$ [5].
3: Compute a residual topology $G^- \leftarrow \{e_G \in G \text{ and } e_G \in G^+\}$ by keeping only edges (i.e., wavelengths) $e_G \in G$
   belonging to output graph $G^+$. Identify an edge set $\mathbb{E} \leftarrow \{e_G \in G \text{ and } e_G \notin G^+\}$ that need to be reconfigured.
4: Reroute traffic on residual topology. Drain input ports connecting $\mathbb{E}$ during the draining latency.
5: Switch wavelengths in $\mathbb{E}$ to update topology $G^+$. Reroute traffic on updated topology.

---

## 2. Topology and Traffic Engineering (TTE) in MAESTRO

DC and HPC topologies can be divided into two groups: one where OCSs interconnect ToRs directly and another where OCSs interconnect aggregation switches. In this paper, we consider a flat reconfigurable Hyper-X network, i.e., a direct connection topology, because it does not require aggregation switches and offers better reconfigurability for skewed inter-rack traffic [1]. The top of Fig. 1(a) shows a schematic of software-defined control and network management functionalities together with data plane network architecture. To perform a hitless reconfiguration, a TTE controller applies traffic and routing configurations based on MAESTRO reconfiguration workflow and computes a new network topology. The historical and real-time traffic and link utilization monitoring can be used to estimate the demand matrix. The TTE controller can use an OpenFlow interface to update routing and forwarding tables on ToR switches for flow rerouting as well as to add or drop wavelengths to reconfigure the network. The Hyper-X network can be built by interconnecting $N$ servers in a rack connected with a $k$ port ToR switch shown in Fig. 1(a). Multiple racks are organized into clusters, where a photonic switch interconnects (black lines) $P$ ToR electrical switches in each row cluster. These row clusters are connected using an additional layer of inter-cluster photonic switches, shown in green links. Thus, these clusters are arranged into $C$ rows and $P$ columns of racks. For example, Fig. 1(b-d) shows a 2D Hyper-X network consisting of six ToRs. The ToRs are indexed with row $i = \{1,2\}$ and column $j = \{1,2,3\}$, and shown in blue circles. Furthermore, they are interconnected through optical switches with links shown in green and black.

Figs. 1(b-d) show TTE operations devised in MAESTRO. Consider a fat flow $f_2$ between $ToR_{1,2}$ and $ToR_{1,3}$ that requires an additional link's bandwidth to satisfy its demand. In this example, we could add a link between the two ToR pairs by breaking two links shown in the dashed black lines in (b). However, note that the traffic on a flow $f_1$ will get disrupted. Therefore, before breaking the links, MAESTRO reroutes flow $f_1$, let these links drain, and adds a red link between $ToR_{1,2}$ and $ToR_{1,3}$. Thus, the software-controlled MAESTRO procedure involves a drain latency in addition to a switching and routing table update latency (which can take from $\mu$s to a few ms), which occurs in any other reconfiguration scheme. Algorithm 1 explains the MAESTRO procedure. In summary, it performs the following operations: i) identify links to be disabled and enabled, ii) reroute flows on a residual topology and drain links to be broken, and iii) reconfigure links and reroute flows on an updated topology. We use an equal-cost multi-path routing for the forwarding of packets, which splits flows over available parallel links between a ToR pair as well as over equal-cost routes between source-destination servers.

## 3. Simulation Setup and Results

We evaluated the MAESTRO and OSR methods using the Netbench packet simulator, and included the routing table update mechanism and reconfiguration evaluation. We consider Flex-LIONS [6] interconnected two different Hyper-X networks with 64 racks arranged into eight rows and eight columns (Fig. 1(b) with $C = P = 8$). Each rack has eight servers interconnected by eight 100Gbps downlink ports of a ToR switch. A Flex-LIONS OCS offers a free-spectral-range (FSR)-based wavelength pairs that can switch over two ports. Thus, we consider a 1FSR Hyper-X network employing commodity 24-port ToR switches with each of 16 uplinks ports having a line rate of $B = 100$Gbps (Fig. 2a, top). Additionally, we consider a 2FSR Hyper-X network with $40-$port ToRs for interconnecting ToR pairs with two parallel links via an OCS, each operating at 50Gbps. While one link maintains fixed connectivity, the other offers reconfigurability. Thus, both networks offer 1.6 Tbps uplink and 0.8 Tbps downlink capacity, i.e., an over-subscription ratio of 2:1. We assumed a link propagation delay of 10ns. We adopted a data center TCP as the transport control protocol between servers with a maximum packet size of 1500B. We set the buffer size of each ToR output port as 10MB, the congestion notification threshold as 100KB for the 1FSR network. For a fair comparison, they are reduced to halves for the 2FSR network. We generated two sets of traffic by deriving the flow size distributions (Fig. 2b) and the overall source-destination traffic pair probability distribution (Fig. 2c-d) from real HPC applications' traces of AMR, and Nekbone [7]. The message's size generated by these applications is emulated by generating flows with approx. the same flow size distribution. However, we assume a Poisson flow arrival process with mean rate $\lambda$. The network load $L$ is given by $\frac{\lambda * F}{S * B}$, where $F, S$ and $B$ are the mean flow size, the number of servers, and the capacity of each server-to-ToR link, respectively.
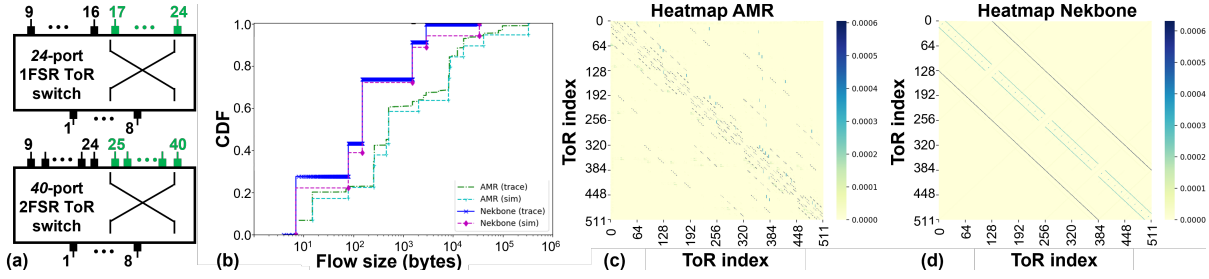
Fig. 2. (a) 1FSR and 2FSR ToR switches schematic. (b) Flow size distribution of AMR, and Nekbone HPC applications. (c-d) Heatmaps of flow-pair distributions of AMR and Nekbone.
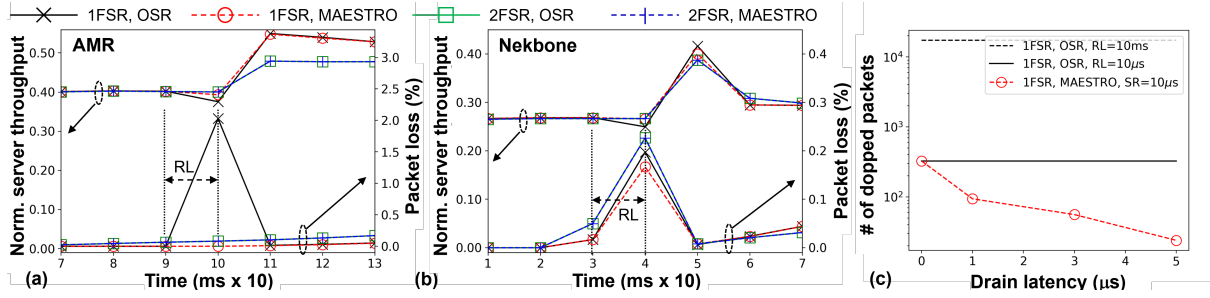


Fig. 3. (a-b) Avg. server throughput is shown when OSR and MAESTRO are performed on 1FSR and 2FSR networks at time=90 (30) ms for a 10ms duration under AMR (Nekbone) traffic with 50 (25)% load. (c) Effect of drain time and RLs on packet loss under AMR traffic at 25% load.

Figs. 3(a-b) show the effect of MAESTRO and OSR in 1FSR and 2FSR networks on the throughput and packet loss based on packets statistics collected every 10 ms. The performance of AMR (and highly skewed Nekbone) is shown for 50(and 25)% load at which they start showing the benefit from reconfiguration in throughput. The reconfiguration latency (RL) is 10ms with no drain latency. After RL, the topology changes to a new pre-computed topology [5]. We make the following observations. i) MAESTRO reduces the throughput dip compared to OSR during the reconfiguration in all scenarios. ii) It keeps the packet loss at the pre-reconfiguration level for AMR and reduces loss due to reconfiguration by $\sim 98\%$ in the 1FSR network. iii) 2FSR network offers a higher resiliency to performance degradation during reconfiguration but offers limited reconfigurability. Nevertheless, a significant change in topology may lead to reduced packet loss improvement by MAESTRO, e.g., Nekbone. To support the highly skewed traffic profile of Nekbone, the post-reconfiguration topology blocks and reconfigures 56% (or 28%) links of prior 1FSR (or 2FSR) Hyper-X topology. Thus, a drastic change in topology should be handled in multiple steps, which might increase RL. Figs. 3(c) illustrates the impact of drain time with various RLs (= switching and routing, i.e., SR + drain latency) on the packet loss at 25% load with AMR traffic. We see the draining effect as the # of dropped packets decreases with increasing RLs. In contrast, we will observe the packet loss increasing due to congestion when RL increases at a higher load.

## 4. Conclusions

We proposed a make-before-break strategy for reconfiguration in optical datacenters (MAESTRO). The simulation results show that MAESTRO may eliminate the excessive packet loss due to reconfiguration by $\sim 98\%$ compared to OSR in 1FSR network. While 2FSR network offers more resiliency to the reconfiguration impact. In the future work, we will study the impact of MAESTRO with applications running on a testbed.

## References

1. M. Y. Teh *et al.*, "Performance trade-offs in reconfigurable networks for hpc," *JOCN*, vol. 14, no. 6, pp. 454–468, 2022.
2. G. Liu *et al.*, "3d-hyper-flex-lion for reconfigurable all-to-all hpc networks," in *SC20*, 2020, pp. 1–16.
3. L. Poutievski *et al.*, "Jupiter evolving ..." in *Proceedings of the ACM SIGCOMM 2022 Conference*, 2022, pp. 66–85.
4. S. K. Singh and A. Jukan, "Machine-learning-based prediction," *JOCN*, vol. 10, no. 10, pp. D12–D28, 2018.
5. S. K. Singh *et al.*, "Multi-cluster reconfiguration with traffic prediction in hyper-flex-lion architecture," in *OFC*, 2022.
6. X. Xiao *et al.*, "Multi-fsr silicon photonic flex-lions," *JLT*, vol. 38, no. 12, pp. 3200–3208, 2020.
7. NERSC, "Characterization of the doe mini-apps. https://portal.nersc.gov/project/CAL/doe-miniapps.htm."