

Predictive Load Balancing for RDMA Traffic

Erfan Nosrati
University of Calgary
Calgary, Canada
erfan.nosrati@ucalgary.ca

Majid Ghaderi
University of Calgary
Calgary, Canada
mghaderi@ucalgary.ca

ABSTRACT

Fast training of large machine learning models requires distributed training on AI clusters consisting of thousands of GPUs. The efficiency of distributed training crucially depends on the efficiency of the network interconnecting GPUs in the cluster. These networks are commonly built using RDMA following a Clos-like datacenter topology. To efficiently utilize the network bandwidth, load balancing is employed to distribute traffic across multiple redundant paths. While there exists numerous techniques for load-balancing in traditional datacenters, these are often either optimized for TCP traffic or require specialized network hardware, thus limiting their utility in AI clusters.

This paper presents the design and evaluation of Hopper, a new load-balancing technique optimized for RDMA traffic in AI clusters. Operating entirely at the host level, Hopper requires no specialized hardware or modifications to network switches. It continuously monitors the current path for congestion and dynamically switches traffic to a less congested path when congestion is detected. Furthermore, it incorporates a lightweight mechanism to identify alternative paths and carefully controls the timing of path switching to prevent excessive out-of-order packets. We evaluated Hopper using ns-3 simulations and a testbed implementation. Our evaluations show that Hopper reduces the average and 99-percentile tail flow completion time by up to 20% and 14%, respectively, compared to state-of-the-art host-based load balancing techniques.

1 INTRODUCTION

Modern machine learning (ML) models, such as large language models (LLMs) and deep learning recommendation models (DLRMs), contain billions and even trillions of parameters. Fast training of these models requires distributed training on AI clusters consisting of thousands of GPUs [8]. A critical challenge in distributed training is the communication overhead during GPU synchronization. In large clusters running multiple training jobs, this communication can dominate the training time [3], significantly slowing the process. The key to faster distributed training is optimizing the cluster network connecting GPU servers to minimize communication overhead.

Large AI clusters are designed following best practices in high-performance computing (HPC) and datacenter networking. Specifically, they adopt Remote Direct Memory Access (RDMA) for high-speed, low-latency communication [8]. Modern AI clusters commonly implement RDMA using RoCEv2 [8] due to its compatibility with standard Ethernet. They also utilize hierarchical datacenter topologies (e.g., leaf-spine [11]) for flexible scale-out [33]. These topologies provide multiple paths between any two GPU servers for fault tolerance, performance, and scalability. To maximize network

bandwidth utilization and thereby reducing inter-GPU communication time in ML training, it is crucial to implement load balancing in the network to efficiently distribute traffic across different paths [24, 33, 37]. While there is extensive work on load balancing in traditional datacenters [1, 9, 15, 34], these works generally consider CPU-generated traffic, such as TCP, and thus do not perform efficiently with GPU-generated RDMA traffic [37]. Currently, there is a gap in the literature on practical and deployable load balancing techniques for RDMA in AI clusters, which we aim to address in this work.

A widely used technique for load balancing in today’s datacenters is Equal Cost Multi-Path (ECMP) [12, 42]. However, numerous studies have shown that ECMP is unable to distribute the load evenly over different paths [34] when flow size distribution is skewed, with a small number of large flows accounting for a disproportionate share of total traffic [10, 15, 34]. This problem is particularly pronounced in distributed ML training workloads, which typically generate fewer but significantly larger flows compared to conventional datacenter workloads, leading to more imbalanced load distribution under ECMP [8, 19].

Several works have been proposed to address the shortcomings of ECMP. For instance, some works use random packet spraying (RPS) [7, 9] to distribute packets uniformly across all paths, achieving near-perfect load balancing, but introduce a significant number of out-of-order (OOO) packets [7, 24, 37]. Other works [1, 38] split a flow into flowlets based on inter-packet time gaps in the flow, and send flowlets over different paths. Although flowlet switching is efficient in balancing the load without excessive OOO packets, its efficiency depends on traffic characteristics, *i.e.*, whether there are flowlets available. This is problematic in AI clusters as hardware-generated RDMA traffic does not include sufficient idle gaps to allow effective flowlet switching [24, 37]. A few works have thus considered switching paths dynamically whenever the current path becomes congested [15, 24, 37]. However, these approaches either choose paths to switch to blindly, irrespective of their quality (FlowBender [15]), necessitate complex host-side logic to manage multiple paths (MP-RDMA [24]), or rely on less widely deployed programmable network switches (ConWeave [37]).

In this work, we present the design and evaluation of Hopper¹, a host-based single-path RDMA load balancing technique. While Hopper cannot match the performance of hardware-based techniques such as ConWeave, it offers simplicity, practicality, and deployability in current AI clusters. Unlike RPS, Hopper is robust to changes in path characteristics and topology asymmetries. Moreover, it outperforms other state-of-the-art solutions like FlowBender [15] (host-based) and CONGA [1] (switch-based). As a host-based solution, Hopper leverages measured path round-trip-time (RTT) as a

¹At the core of our technique is *hopping* over congested network paths, hence the name Hopper.

signal to detect congested paths and reroute flows to less utilized alternative paths, while minimizing the number of OOO packets during path switching. Hopper requires only standard ECMP support on network switches and minimal host modifications for RTT estimation and path switching. Both of these functionalities have been well-studied by prior work [15, 18, 19, 25, 28, 34, 35], with a number of efficient solutions already available.

Our work is motivated by FlowBender [15], which targets TCP traffic. Similar to Hopper, it performs path switching based on congestion signals (*i.e.*, ECN marking), but employs a random path selection strategy, which we argue can lead to degraded performance for two key reasons:

- **Suboptimal Path Selection:** The newly selected path might also be congested. In such cases, FlowBender continues to use the new congested path for a few RTTs before it decides to switch again, prolonging performance issues.
- **Significant Out-of-Order Packets:** There is a potential for a large difference in path delays between the old and new paths. Thus, with high transmission rates common in RDMA, switching to a random path can cause a significant number of OOO packets.

To avoid these issues, Hopper employs an *informed* path selection and switching strategy by leveraging the capabilities of modern RDMA NICs (RNICs) for RTT measurement and limited OOO packet handling. First, it only initiates path changes when a clearly better alternative is available, thereby avoiding unnecessary rerouting under heavy network load when all paths are similarly congested. Specifically, when the current path RTT exceeds a certain threshold, Hopper begins probing for less congested paths by sending small probe packets on two randomly selected alternative paths. Once a path is identified as congested, *i.e.*, its RTT exceeds a second threshold, Hopper compares the congestion level of the current path with that of two recently probed alternatives. If either probed path offers significantly better conditions, Hopper changes the path of the flow to the best of the two. Otherwise, the flow continues on the current path. Second, Hopper carefully considers the delay difference between the old and new paths. Specifically, packets on the new path are delayed in proportion to the RTT difference between the current and new path, reducing the chance of OOO packets that could trigger retransmissions. Using simulations and testbed experiments, we show that these design optimizations lead to substantial performance gains without introducing any significant overhead, enabled by the capabilities of modern RNICs for packet reordering and RTT measurement.

The contributions of this paper are summarized as follows:

- We design and evaluate Hopper, a host-based load balancing technique for AI clusters that operates at near-RTT granularity, without causing excessive OOO packets.
- We develop mechanism for light-weight congestion-aware path selection and switching by leveraging the capabilities of RNICs for RTT measurement and limited packet reordering.
- We evaluate Hopper using both ns-3 simulations and a hardware testbed implementation. Our results show that, compared to FlowBender, Hopper improves average and 99-percentile flow completion time (FCT) by up to 20% and 14%, respectively, for ML training workloads.

The rest of the paper is organized as follows. We discuss our motivation for Hopper in §2. Design principles and components of Hopper are discussed in §3. Evaluations are presented in §4. Related works are reviewed in §5, while §6 concludes the paper.

2 BACKGROUND AND MOTIVATION

The primary motivation for our work stems from the observation that random path selection as employed in RPS [7], FlowBender [15], and PLB [34], while simple, is not efficient. We argue that with recent advances in RNICs, such as hardware time-stamping and packet reordering, it is now feasible to efficiently implement congestion-aware path selection on the hosts without requiring any support from network switches, modifying packet headers, or degrading the NIC's performance. In this section, we provide a brief overview of RDMA and discuss our motivation in more detail.

RDMA. RDMA is a high-performance communication technology that delivers low latency and high throughput by implementing transport functionalities including congestion control and loss recovery entirely in NIC hardware. It enables hardware-accelerated direct memory access over the network without involving the CPU or operating system. An RDMA connection is identified by a so-called Queue Pair (QP). The QP is the communication interface that allows user-space applications to send and receive messages via the RDMA transport residing on the RNIC hardware.

RoCEv2. RDMA over Converged Ethernet version 2 (RoCEv2) encapsulates RDMA messages within UDP/IP packets, enabling RDMA to operate over standard Ethernet-based datacenter networks. RDMA was originally designed to operate over lossless Infiniband fabrics [2]. To enable RDMA over Ethernet, RoCEv2 relies on Priority Flow Control (PFC) [13] to turn Ethernet into a lossless fabric. However, PFC brings challenges such as congestion spreading and deadlocks [26]. To address these issues, several mechanisms are added to RDMA network stack including end-to-end congestion control like DCQCN [46], and selective repeat loss recovery like IRN [26].

Problem with Out-of-Order Packets. Modern RNICs support packet reordering in the hardware via IRN. However, to avoid degrading the NIC's performance, it is crucial to limit the number of OOO packets. IRN requires per-connection memory for tracking lost packets and rearranging OOO packets. Commercial RNICs have limited on-chip SRAM memory [40]. For example, NVIDIA CX-5 RNIC has only about 2 MB on-chip memory [16]. Allocating even one BDP (bandwidth-delay product)-sized reordering buffer consumes 0.1 MB memory (network bandwidth 100 Gbps and RTT 8 μ s) per connection, severely limiting the number of connections the RNIC can support.

Problem with Random Packet Spraying. RPS looks like the perfect load balancing approach as it uniformly distributes the load across all paths. But, that is exactly where the problem lies. Not only this creates a massive number of OOO packets [37], overwhelming the RNIC resources, but also is vulnerable to network asymmetries. For uniform load balancing, RPS requires a symmetric network topology with consistent path characteristics. However, in large AI clusters, device failures (*e.g.*, links and switches) are common [33], leading to an asymmetric topology. Moreover, transient congestion due to rerouted traffic can also lead to asymmetries in path

Alg. 1 Hopper Control Logic in Each RTT Epoch.

```

avg_rtt ← 0, probe ← true, switch ← true
for every new_rtt measurement do
  avg_rtt ←  $\alpha \cdot \text{new\_rtt} + (1 - \alpha) \cdot \text{avg\_rtt}$ 
  if avg_rtt > th_probe & probe == true then
    alt_paths ← PROBEPATHS()
    probe ← false
  end if
  if avg_rtt > th_cong & switch == true then
    SWITCHPATH(alt_paths)
    switch ← false
  end if
end for

```

characteristics. With such asymmetries, uniformly sending packets over all paths leads to inflated tails for flow completion times. This degradation is particularly detrimental for distributed ML training, where training progress is gated by the completion time of the slowest flow in a collective communication operation [5].

Summary. An ideal solution is to make RPS consider path congestion on a packet-by-packet basis. But the overhead of such an approach is prohibitive. A more practical solution is to increase the load balancing granularity to amortize the path assessment overhead. This is exactly what Hopper aims to achieve by operating at RTT granularity, slower than per-packet but faster than per-flow load balancing. While Hopper can't achieve perfect uniform load distribution as RPS, it does not create as many OOO packets and is robust against topology changes.

3 DESIGN

Our design for Hopper leverages the ability of modern RNICs [29, 30] for limited packet reordering and precise time-stamping. In the following subsections, we present the design of Hopper and discuss its various components. An example of Hopper's workflow is presented in Appendix § A.

Design Overview. Hopper operates over control epochs, where the duration of each epoch is set to one RTT. It consists of three main modules: (1) a congestion detection module, (2) a path probing module, and (3) a path switching module. Together, these modules enable Hopper to dynamically reroute flows to less congested paths by only manipulating packet headers on the host to affect standard ECMP hashing decisions. The high-level pseudo-code for Hopper is presented in Alg. 1. Below, we describe the functionality of each module.

3.1 Congestion Detection Module

Congestion Signal. The two most widely used congestion signals are Explicit Congestion Notification (ECN) and Round-Trip-Time (RTT). These signals are often readily available at RDMA transport layer, but could also be obtained independently with low overhead [35]. For example, DCQCN [46] relies on ECN marking by switches, while Timely [25] uses RTT measurements to infer congestion. Hopper leverages RTT for path-level congestion detection, as it is supported by commodity NICs [18], requires no switch support, and is particularly suited for path probing. Since probing

typically involves transmitting only a small number of packets, ECN may not reliably indicate congestion due to insufficient marking opportunities. On the other hand, ECN is capable of signaling emerging congestion earlier, even before RTT increases become apparent [19]. While Hopper can be further optimized to leverage ECN markings to initiate path probing proactively, our current implementation is based on RTT measurements.

Congestion Detection. Hopper monitors measured per-packet RTT (denoted by new_rtt in Alg. 1) on the sending host, *e.g.*, using ACKs in IRN [26]. It keeps a moving average of measured RTTs over a control epoch of one RTT. The current path is considered congested when the average RTT exceeds a predefined threshold th_cong. The threshold th_cong can be tuned to best match the specific network architecture and workload characteristics. Also, the moving average parameter α can be increased (decreased) to make Hopper more (less) responsive to RTT increases.

3.2 Path Probing Module

Probe Mechanism. Given modern RNICs' support for limited out-of-order delivery, probing can be integrated directly into data transmission, sending a few data packets along alternate paths. Therefore, in our current design, Hopper uses out-of-band probing packets to scan alternative paths. However, sending probing packets, if not carefully controlled, can increase network load, especially when the network is already congested. Moreover, probing a new path in RDMA requires using a new QP, which could affect the scalability of RNICs [40], if not controlled carefully.

Probe Initiation. To address these challenges, Hopper employs the power-of-two-choices strategy. If the average RTT exceeds a threshold th_probe, probing for alternate paths starts. To this end, Hopper probes two randomly selected unexplored paths by creating two QPs bound to distinct UDP source ports. It records these source ports along with the corresponding measured delays to inform rerouting decisions. In our testbed implementation, we perform profiling to map specific source ports to distinct paths, ensuring consistent path selection during probing. To prevent redundant probing, Hopper employs timestamping to track recently explored paths and avoids selecting any path again if it was probed within the last ttl_probe interval.

3.3 Path Switching Module

Path Selection. Hopper switches the current path only if it is congested and there is a less congested path available. Otherwise, it will wait for one RTT epoch to initiate probing again. To prevent rerouting to paths exhibiting similar congestion levels, Hopper initiates rerouting only when the alternative path demonstrates a substantially lower RTT than the current path, exceeding a configurable margin δ_{rtt} . This parameter can be optimized to balance performance and stability by avoiding path switches that do not result in meaningful improvements. By preventing unnecessary path switches, this approach would also reduce the number of out-of-order packets.

Path Switching. To switch to a new path, Hopper simply starts using the corresponding QP that was used to probe the path and releases the QP of the old path. If neither path offers improvement, Hopper frees the probed QPs but retains the source port

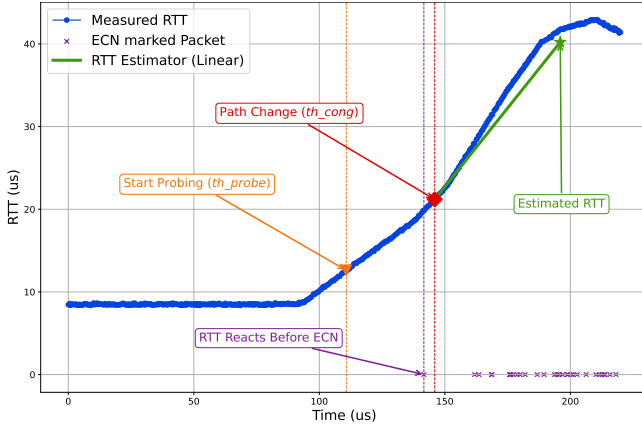


Figure 1: Hopper employs a linear estimator to predict path RTT.

and delay information for a few RTTs to avoid repeatedly probing the same congested paths. Although this approach may introduce some delay in switching away from congested paths, it effectively prevents excessive path switching when all links are persistently congested. Additionally, it avoids the issue where multiple flows from congested paths are rerouted to the same alternative path due to aggressive probing frequency.

Reducing Out-of-Order Packets. Since the RTTs of both the current and alternative paths are available, Hopper can delay rerouting by a duration proportional to the RTT difference between the two paths. This delay helps reduce out-of-order packet delivery at the receiver [24]. However, if not carefully implemented, such delayed rerouting can introduce unnecessary latency, contradicting RDMA's low-latency goals and degrading end-to-end performance. To calculate a suitable delay, Hopper needs to estimate RTT on the current path for the in-flight packets. As depicted in Fig. 1, we employ a linear regression model for this purpose. The model calculates the rate of RTT increase using `new_rtt` measurements in the current RTT epoch and extrapolates this trend based on the number of in-flight packets to estimate the RTT of the last packet on the congested path. Although the estimated RTT may overestimate the actual path delay, since RTT increases tend to stabilize once queues stop growing, Hopper uses this value as a conservative upper bound.

4 EVALUATION

We use ns-3 [27] simulations and a tested implementation to evaluate Hopper under traditional datacenter and ML workloads. The objective of our evaluation is to address the following questions:

- *How much improvement congestion-aware path selection can achieve compared to random path selection?* To answer this question, we compare Hopper with FlowBender [15].
- *How effective dynamic load-balancing is compared to flowlet switching for RDMA?* To answer this question, we compare Hopper with CONGA [1].
- *What is the performance gap between host-based and in-network load balancing?* To answer this question, we compare with ConWeave [37].

Table 1: Hopper parameters.

Parameter	Value
α	1
th_prob	$1.5 \times \text{base RTT}$
th_cong	$2.5 \times \text{base RTT}$
ttl_prob	$4.0 \times \text{base RTT}$
δ_{rtt}	80%

4.1 Simulation Experiments

We build on the ns-3 implementation provided by ConWeave, using the same network topology, transport parameters, and workload generator. However, in addition to datacenter workloads used in ConWeave, we also incorporate an ML training workload to specifically study the behavior of Hopper and other techniques in AI clusters. The parameters used for Hopper in the simulations are shown in Table 1. These values are found to work well in our setup.

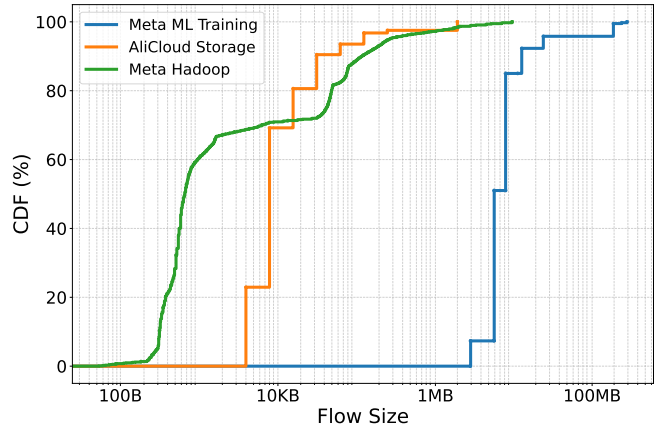


Figure 2: Traffic characteristics of representative workloads. The flow sizes in ML Training workload are substantially larger than those in datacenter workloads.

4.1.1 Setup and Methodology.

Network Topology. We adopt a leaf-spine topology comprising 128 servers and 16 switches. Each group of 16 servers connects to a dedicated leaf switch, and all 8 leaf switches are fully connected to 8 spine switches. All links are 100 Gbps with a latency of $1 \mu\text{s}$, for a base RTT of $8 \mu\text{s}$.

Transport. We use DCQCN [46] for congestion control, which matches the algorithm used by NVIDIA CX-5 RNICs in our testbed. We set DCQCN parameters identical to those used in ConWeave. We also implement and extended version of IRN [26] for packet re-ordering on RNICs. When an OOO packet is received, it is buffered and acknowledged normally as long as its sequence number falls within a predefined threshold (set to 30 packets) relative to the expected sequence number. If the sequence number exceeds this threshold, the receiver initiates loss recovery by sending a NACK containing both a cumulative acknowledgment and a selective acknowledgment (SACK). All other aspects of IRN, including limiting in-flight data to one BDP, remain unchanged.

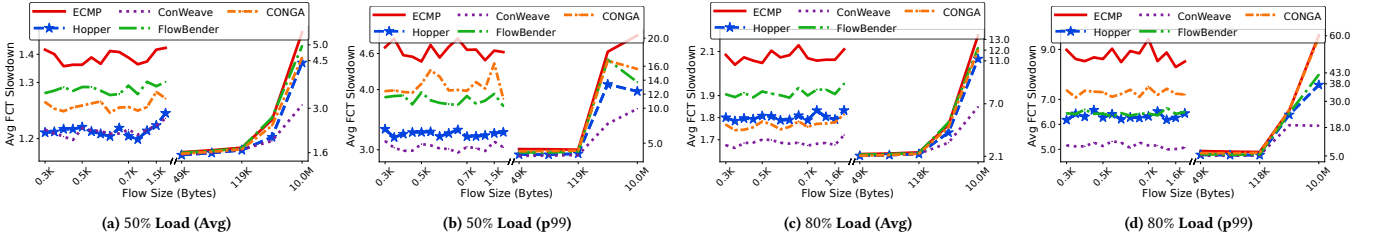


Figure 3: Average and tail FCT slowdown for Hadoop workload at 50% and 80% average network load.

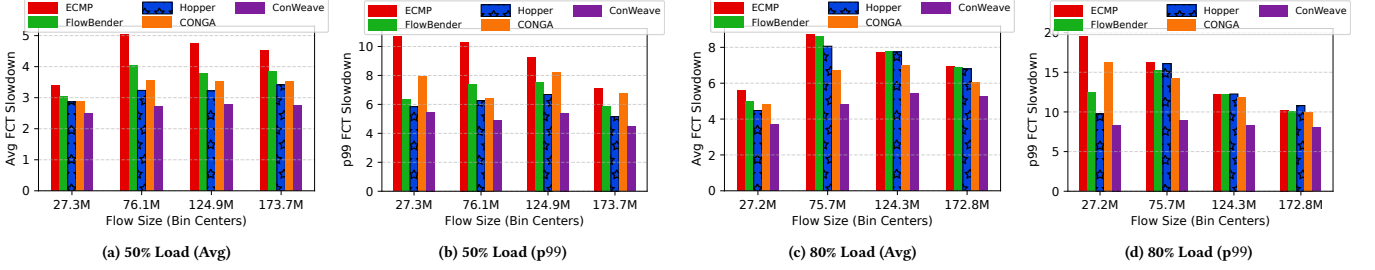


Figure 4: Average and tail FCT slowdown for ML Training workload at 50% and 80% average network load.

Workloads. We conduct experiments using traditional datacenter workloads, including AliCloud Storage [21] and Meta Hadoop [36], as well as traffic patterns representing those of distributed ML training. Specifically, we utilize the collective message size distribution reported in [8], which analyzes ML training jobs of up to 128 GPUs on Meta’s internal AI cluster. The workload includes a mix of collective communication operations [6], with AllReduce being common in DDP [20], while AllGather and ReduceScatter being prevalent in FSDP [45]. Message sizes vary depending on the model architecture and training strategy, capturing a broad range of realistic communication demands. The CDFs of these workloads are shown in Fig. 2. We used the ConWeave’s load generator to simulate network traffic by randomly selecting a sender–receiver pair and assigning a flow size based on the workload distribution. Flow start times follow a Poisson distribution and are chosen to create two network load scenarios, a moderate load scenario at 50% and a high load scenario at 80% network load.

Baseline Schemes. We compare Hopper with FlowBender, CONGA, and ConWeave. We do not include PLB in our evaluation, as it is specifically designed for TCP flows and behaves similar to ECMP in RDMA. Also, like FlowBender, it employs a random path selection strategy.

Performance Metric. We use flow completion time (FCT) slowdown as the primary performance metric, defined as the ratio of a flow’s actual completion time to its baseline completion time in an unloaded network. We present both the average and 99-percentile tail (p99) of the slowdown.

4.1.2 Results and Discussion.

Datacenter Workloads. Fig. 3 shows the results for Hadoop workload. The results for AliCloud workload are similar and included in the Appendix B. From the workload distribution in Fig. 2, we observe that only a small percentage of flows are between 2KB and 49KB. Thus, for the sake of presentation, the x-axis in the plots is broken into two regions, the left region shows short flows (< 2KB),

while the right region shows large flows (> 49KB). The largest flow size in this workload is 20MB, and only less than 5% of flows are larger than 266KB. Therefore, we caution that the reported results for 10MB flow size are not stable. Instead, for large flow sizes, we focus on the results reported for the ML Training workload in the next subsection. The main observation is that Hopper consistently outperforms FlowBender for both average (by up to 7.8%) and p99 tail slowdown (by up to 19.6%) across different loads and flow sizes. Interestingly, even though both Hopper and FlowBender operate at RTT granularity, and thus do not load balance short flows, we see significant improvement compared to ECMP due to them load balancing large flows more effectively. We also observe that Hopper similarly outperforms CONGA except for the average slowdown in high load, where they are both comparable. Finally, in moderate load, Hopper’s performance is comparable to ConWeave for small flows, with a gap of less than 8.3%.

ML Workload. For the ML Training workload, traffic is more concentrated around specific flow sizes. Thus, for presentation clarity, as depicted in Fig. 4, we divide flows into four bins to show FCT slowdown per size category. Recall that this workload mainly consists of large flows that represent collective communication operations. Therefore, not surprisingly, the behavior of different approaches is consistent with the scenario of large flows in datacenter workloads. Specifically, we observe that 1) Hopper outperforms FlowBender across all scenarios, achieving up to 20% and 14% improvement in average and p99 FCT slowdown, respectively, and 2) when the network is moderately loaded, Hopper even outperforms CONGA by up to 10% and 23% in average and p99 FCT slowdown, respectively. However, with high network load, switch-based techniques generally outperform host-based ones, specially for the larger flow sizes.

4.2 Testbed Experiments

We implement a software-based prototype of Hopper on a hardware testbed, as depicted in Fig. 5. The hardware components of

the testbed include Dell PowerEdge R740 servers equipped with NVIDIA CX-5 RNICs [29], and a Dell PowerSwitch S5248F-ON switch running SONiC OS [32].

4.2.1 Setup and Methodology.

Implementation. Due to the black-box nature of CX-5 RNICs, implementing Hopper directly within the RNIC hardware is not feasible. While some prior works [23, 39, 44] implement similar mechanisms using software frameworks such as DPDK [31], we choose to preserve compatibility with the RDMA communication model. To that end, the control logic of Hopper is implemented at user space leveraging the Linux rdma-core library [22], while the actual data transfers are implemented using standard RDMA verbs. To send a flow, we divide it into chunks and send chunks via RDMA. The chunk size specifies the granularity at which Hopper can switch paths. We report experiment results with two chunk sizes for data transfers, namely 10MB and 1MB. We find that chunk sizes less than 100KB degrade performance due to increased frequency of system calls and context switches. For probing alternative paths, we use a 10KB chunk size to reduce probing overhead. To measure RTT on a path, we measure the time between initiating an RDMA send operation and receiving the corresponding completion event in the completion queue.

Network Topology. Our physical testbed has a leaf-spine topology (see Fig. 5) with two leaf switches and six spine switches, created on our physical switch using VRF. Each RNIC is equipped with two 25 Gbps ports, each assigned to one host. Each leaf is connected to four hosts in a rack. Also, each leaf is connected to four spine switches via 10 Gbps links and to the remaining two spine switches via 1 Gbps links, creating deliberate asymmetry in path bandwidths.

Transport. All experiments use DCQCN as the congestion control algorithm, configured with the default parameters recommended by NVIDIA for CX-5 RNICs [29].

Workload. We evaluate Hopper using an AllReduce collective communication pattern, where each sender transmits data to its corresponding receiver. The size and frequency of these collective operations are derived from the GPT-3 model training traces [41]. We extract the communication sizes from the traces and generate flow workloads using a custom script. The workload consists of 204 flows, and we repeat each experiment five times, reporting the average to mitigate the effects of randomness. ML training consists of several rounds. In each round, nodes synchronize their intermediate computations using AllReduce before starting a new round. In our testbed, to simplify implementation, at the end of each round of training, clients send a completion message to a centralized server. The next round begins only after the server confirms that all clients have completed the previous round, *i.e.*, have completed their associated flows.

Baseline Schemes. We compare Hopper against two baselines, namely ECMP and a simplified software version of FlowBender. In our FlowBender implementation, congestion is inferred from RTT increases instead of ECN markings. As mentioned earlier, ConWeave and CONGA both require switch modifications, ConWeave in particular requires P4 switches, which we do not have in our testbed. Thus, our testbed comparison is limited to ECMP and FlowBender.

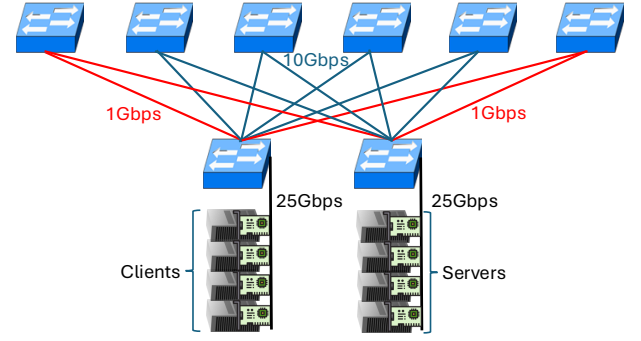


Figure 5: The leaf-spine topology of the testbed with asymmetric path bandwidths.

Performance Metrics. We use the following performance metrics for comparison: 1) the average and tail FCT, 2) the average total training time, and 3) the average link utilization. Specifically, we examine the utilization of 1G links, as their use indicates poor load balancing. Since transmitting flows from all four clients over a single 10G link offers higher throughput than on 1G links, higher use of 1G links reflects inefficient path selection by the load balancing technique.

4.2.2 Results and Discussion.

Avoiding Congested Paths. When a flow starts, it is initially assigned to a random path, which potentially could be a 1G path. However, both FlowBender and Hopper will try to switch the flow to other paths after a short time. Figs. 6(a) and 6(c) show the average utilization of 1G and 10G links in the network when using different chunk sizes. We make the following observations. First, Hopper consistently utilizes 10G links more than FlowBender, by up to 35%, clearly demonstrating its ability to discover less congested paths. Second, as the chunk size increases so does the 1G link utilization under both Hopper and FlowBender. This is because whenever these algorithms choose a 1G path, they stay on that path for a longer time when using large chunks. Yet, the performance degradation of FlowBender is more pronounced, almost twice that of Hopper, due to it having a higher chance of choosing 1G paths, as it employs a random path selection strategy.

FCT Slowdown. As shown in Figs. 6(b) and 6(d), Hopper outperforms FlowBender in all scenarios for both average and tail slowdown, except for p99 with 10MB chunks, where it matches FlowBender. As with link utilization, the penalty of choosing a congested path with 10MB chunk is higher for both approaches, thus leading to a smaller gap between their performances. We note that, smaller chunk sizes result in more granular load balancing decisions, but also higher system overhead in our implementation. For example, with ECMP, we see from the figures that its slowdown increases for 1MB chunks, as it does not benefit from smaller chunk sizes. For both Hopper and FlowBender, the chunk size provides a trade-off between load balancing granularity and system overhead. In our testbed, 1MB chunk size provides a good trade-off. Specifically, with 1MB chunks, Hopper improves the average, p95, and p99 FCT slowdown by 45%, 61.9%, and 77.2%, respectively, compared to FlowBender.

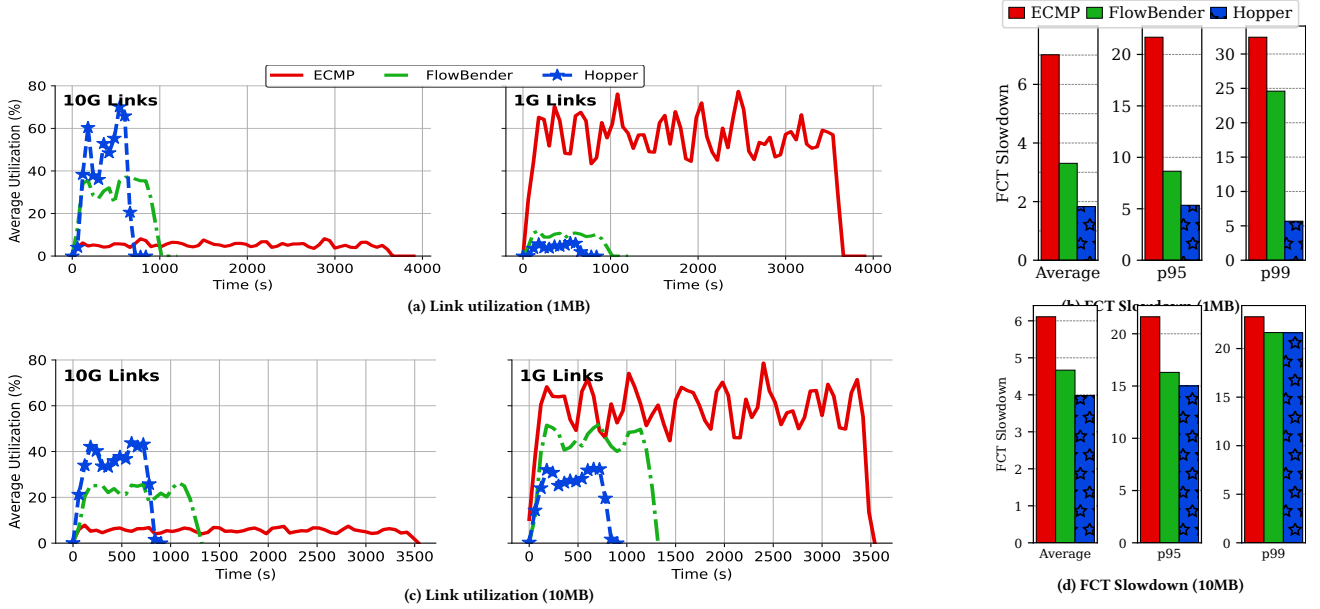


Figure 6: Performance of different load balancing approaches in the testbed with 1MB and 10MB chunks.

Training Time. Figs. 6(a) and 6(c) also show the average total training time (multiple training rounds) under different load balancers. Specifically, with 1MB chunks, the average training time is 1080 and 715 seconds under Hopper and FlowBender, representing 51% reduction in training time. Even with 10MB chunks, the average training time is 1320 and 890 seconds under Hopper and FlowBender, representing 48.3% reduction in training time. This clearly demonstrates the superiority of Hopper for ML training workloads.

5 RELATED WORK

In addition to the works already mentioned in the Introduction [1, 7, 9, 15, 24, 34, 37, 38], we briefly review a few other works that are related to our work on Hopper.

Flowlets in RDMA. As discussed earlier, RDMA traffic lacks sufficient inter-packet gaps to identify flowlets. This observation is confirmed in several works using real measurements [24, 37]. So, one idea then is to artificially insert such gaps in the packet stream in order to create flowlets for load balancing, as recently proposed in HF²T [4]. However, adding flowlet delays, specially when the network is not congested, could unnecessarily prolong FCTs and still relies on the existence of a flowlet load balancing, such as CONGA [1], in the network.

Multipath RDMA. MP-RDMA [24] includes several sophisticated mechanisms to break a flow to subflows that are then distributed over multiple paths with respect to the congestion level on those paths. However, it requires complex host-side logic to manage multiple paths and synchronize packet transmissions across them. NVIDIA collective communication library (NCCL) [6], which is widely used for distributed training in large AI clusters, internally uses multiple QPs to send each message [8], but employs random

spraying similar to RPS, albeit at a chunk instead of packet granularity. Thus, in asymmetric networks, it faces similar issues as those faced by RPS.

Path Probing. Both host and switch-based path probing have been extensively considered for various network management tasks. For example, both CONGA [1] and ConWeave [37] implement passive and active path probing on network switches, respectively, for load balancing. On the other hand, Clove [17], MP-RDMA [24], and Hermes [43] implement path probing on the host. For instance, Clove pre-computes a set of disjoint paths before transmission (similar to our testbed implementation) and periodically probes them, while Hermes leverages the power-of-two-choices technique to probe alternative paths with minimal overhead.

RTT Measurement. The hardware timestamping capability of modern NICs [14, 29, 30] has enabled precise measurement of network RTT. Timely [25] demonstrates that, by leveraging this capability, RTT can be measured with high precision. Swift [18] builds on this by using NIC-based timestamping to eliminate the effects of local NIC delays, resulting in more accurate RTT measurements. NVIDIA's ZTR-RTT [28] congestion control also relies on hardware timestamping to accurately measure RTT and dynamically adjust transmission rates.

6 CONCLUSION

This paper presents Hopper, a host-based load balancing technique for RDMA traffic in AI clusters. Hopper operates at RTT granularity, meaning it can switch flow paths every RTT. However, to avoid unnecessarily rerouting flows, it switches paths only when the current path of a flow is congested and a less congested alternative path can be found via a proactive light-weight probing mechanism. Using both simulations and hardware testbed implementation, we evaluated Hopper under traditional datacenter as well as ML training

workloads. Our evaluation revealed that Hopper outperforms state-of-the-art host-based load balancing techniques, while achieving performance close to that of switch-based techniques.

A major challenge in our work was the black-box nature of commercial RDMA NICs. To overcome this, we implemented Hopper's control logic in user space. A potential direction for future research is to extend our user space implementation by incorporating Hopper in a collective communication library, as inter-GPU communication in distributed ML training is entirely based on collective operations regardless of the model architecture or training framework.

REFERENCES

- [1] Mohammad Alizadeh, Tom Edsall, Sarang Dharmapurikar, Ramanan Vaidyanathan, Kevin Chu, Andy Fingerhut, Vinh The Lam, Francis Matus, Rong Pan, Navindra Yadav, and George Varghese. 2014. CONGA: distributed congestion-aware load balancing for datacenters. In *Proceedings of the 2014 ACM Conference on SIGCOMM* (Chicago, Illinois, USA) (SIGCOMM '14). Association for Computing Machinery, New York, NY, USA, 503–514.
- [2] InfiniBand Trade Association. 2007. Infiniband Architecture Specifications. <https://www.infinibandta.org/ibta-specification/>.
- [3] Jiamin Cao, Yu Guan, Kun Qian, Jiaqi Gao, Wencong Xiao, Jianbo Dong, Binzhang Fu, Dennis Cai, and Ennan Zhai. 2024. Crux: GPU-Efficient Communication Scheduling for Deep Learning Training. In *Proceedings of the Conference of the ACM Special Interest Group on Data Communication* (Sydney, NSW, Australia) (ACM SIGCOMM '24). Association for Computing Machinery, New York, NY, USA, 1–15.
- [4] Chuhao Chen, Jiarui Ye, Yongbo Gao, Sen Liu, and Yang Xu. 2024. HF²T: Host-Based Flowlet Fine-Tuning for RDMA Load Balancing. In *Proceedings of the 8th Asia-Pacific Workshop on Networking* (Sydney, Australia) (APNet '24). Association for Computing Machinery, New York, NY, USA, 9–15.
- [5] Xiaoqi Chen, Shay Vargaftik, and Ran Ben Basat. 2024. When ML Training Cuts Through Congestion: Just-in-Time Gradient Compression via Packet Trimming. In *Proceedings of the 23rd ACM Workshop on Hot Topics in Networks* (Irvine, CA, USA) (HotNets '24). Association for Computing Machinery, New York, NY, USA, 177–185.
- [6] NVIDIA Corporation. 2025. NVIDIA Collective Communications Library (NCCL). <https://developer.nvidia.com/nccl>. Accessed: 2025-05-29.
- [7] Advait Dixit, Pawan Prakash, Y. Charlie Hu, and Ramana Rao Kompella. 2013. On the impact of packet spraying in data center networks. In *2013 Proceedings IEEE INFOCOM*. 2130–2138.
- [8] Adithya Gangidi, Rui Miao, Shengbao Zheng, Sai Jayesh Bondu, Guilherme Goes, Hany Morsy, Rohit Puri, Mohammad Riftadi, Ashmitha Jeevaraj Shetty, Jingyi Yang, Shuaiyang Zhang, Mikel Jimenez Fernandez, Shashidhar Gandham, and Hongyi Zeng. 2024. RDMA over Ethernet for Distributed Training at Meta Scale. In *Proceedings of the ACM SIGCOMM 2024 Conference* (Sydney, NSW, Australia) (ACM SIGCOMM '24). Association for Computing Machinery, New York, NY, USA, 57–70.
- [9] Soudeh Ghorbani, Zibin Yang, P. Brighten Godfrey, Yashar Ganjali, and Amin Firoyzshahian. 2017. DRILL: Micro Load Balancing for Low-latency Data Center Networks. In *Proceedings of the Conference of the ACM Special Interest Group on Data Communication* (Los Angeles, CA, USA) (SIGCOMM '17). Association for Computing Machinery, New York, NY, USA, 225–238.
- [10] Junzhi Gong, Tong Yang, Haowei Zhang, Hao Li, Steve Uhlig, Shigang Chen, Lorna Uden, and Xiaoming Li. 2018. HeavyKeeper: An Accurate Algorithm for Finding Top-k Elephant Flows. In *2018 USENIX Annual Technical Conference (USENIX ATC 18)*. USENIX Association, Boston, MA, 909–921.
- [11] Albert Greenberg, James R. Hamilton, Navendu Jain, Srikanth Kandula, Changhoon Kim, Parantap Lahiri, David A. Maltz, Parveen Patel, and Sudipta Sengupta. 2009. VL2: a scalable and flexible data center network. In *Proceedings of the ACM SIGCOMM 2009 Conference on Data Communication* (Barcelona, Spain) (SIGCOMM '09). Association for Computing Machinery, New York, NY, USA, 51–62.
- [12] Christian Hopps. 2000. Analysis of an Equal-Cost Multi-Path Algorithm. RFC 2992.
- [13] IEEE. 2008. 802.1Qbb – Priority-based Flow Control. <https://1.ieee802.org/dcb/802-1qbb/>.
- [14] Intel Corporation. 2025. Intel® Ethernet Network Adapter E810-2CQDA2 Specifications. <https://www.intel.com/content/www/us/en/products/sku/210969/intel-ethernet-network-adapter-e8102cqda2/specifications.html>. Accessed: 2025-06-02.
- [15] Abdul Kabbani, Balajee Vamanan, Jahangir Hasan, and Fabien Duchene. 2014. FlowBender: Flow-level Adaptive Routing for Improved Latency and Throughput in Datacenter Networks. In *Proceedings of the 10th ACM International on Conference on Emerging Networking Experiments and Technologies* (Sydney, Australia) (CoNEXT '14). Association for Computing Machinery, New York, NY, USA, 149–160.
- [16] Anuj Kalia, Michael Kaminsky, and David Andersen. 2019. Datacenter RPCs can be General and Fast. In *16th USENIX Symposium on Networked Systems Design and Implementation (NSDI 19)*. USENIX Association, Boston, MA, 1–16.
- [17] Naga Katta, Aditi Ghag, Mukesh Hira, Isaac Keslassy, Aran Bergman, Changhoon Kim, and Jennifer Rexford. 2017. Clove: Congestion-Aware Load Balancing at the Virtual Edge. In *Proceedings of the 13th International Conference on Emerging Networking Experiments and Technologies* (Incheon, Republic of Korea) (CoNEXT '17). Association for Computing Machinery, New York, NY, USA, 323–335. <https://doi.org/10.1145/3143361.3143401>
- [18] Gautam Kumar, Nandita Dukkipati, Keon Jang, Hassan M. G. Wassel, Xian Wu, Behnam Montazeri, Yaogong Wang, Kevin Springborn, Christopher Alfeld, Michael Ryan, David Wetherall, and Amin Vahdat. 2020. Swift: Delay is Simple and Effective for Congestion Control in the Datacenter. In *Proceedings of the Annual Conference of the ACM Special Interest Group on Data Communication on the Applications, Technologies, Architectures, and Protocols for Computer Communication* (Virtual Event, USA) (SIGCOMM '20). Association for Computing Machinery, New York, NY, USA, 514–528.
- [19] Yanfang Le, Rong Pan, Peter Newman, Jeremias Blendin, Abdul Kabbani, Vipin Jain, Raghava Sivaramu, and Francis Matus. 2024. STrack: A Reliable Multipath Transport for AI/ML Clusters. [arXiv:2407.15266](https://arxiv.org/abs/2407.15266)
- [20] Shen Li, Yanli Zhao, Rohan Varma, Omkar Salpekar, Pieter Noordhuis, Teng Li, Adam Paszke, Jeff Smith, Brian Vaughan, Pritam Damania, and Soumith Chintala. 2020. PyTorch Distributed: Experiences on Accelerating Data Parallel Training. [arXiv:2006.15704](https://arxiv.org/abs/2006.15704) [cs.DC] <https://arxiv.org/abs/2006.15704>
- [21] Yuliang Li, Rui Miao, Hongqiang Harry Liu, Yan Zhuang, Fei Feng, Lingbo Tang, Zheng Cao, Ming Zhang, Frank Kelly, Mohammad Alizadeh, and Minlan Yu. 2019. HPCC: high precision congestion control. In *Proceedings of the ACM Special Interest Group on Data Communication* (Beijing, China) (SIGCOMM '19). Association for Computing Machinery, New York, NY, USA, 44–58.
- [22] Linux RDMA. 2025. RDMA core userspace libraries and daemons. <https://github.com/linux-rdma/rdma-core>. Accessed: 2025-05-16.
- [23] Minfei Long, Jiangping Han, Wentao Wang, Jiayu Yang, and Kaiping Xue. 2024. LSCC: Link-Segmented Congestion Control for RDMA in Cross-Datacenter Networks. In *2024 IEEE/ACM 32nd International Symposium on Quality of Service (IWQoS)*. 1–10. <https://doi.org/10.1109/IWQoS61813.2024.10682909>
- [24] Yuanwei Lu, Guo Chen, Bojie Li, Kun Tan, Yongqiang Xiong, Peng Cheng, Jiansong Zhang, Enhong Chen, and Thomas Moscibroda. 2018. Multi-Path Transport for RDMA in Datacenters. In *15th USENIX Symposium on Networked Systems Design and Implementation (NSDI 18)*. USENIX Association, Renton, WA, 357–371.
- [25] Radhika Mittal, Vinh The Lam, Nandita Dukkipati, Emily Blem, Hassan Wassel, Monia Ghobadi, Amin Vahdat, Yaogong Wang, David Wetherall, and David Zats. 2015. TIMELY: RTT-based Congestion Control for the Datacenter. In *Proceedings of the 2015 ACM Conference on Special Interest Group on Data Communication* (London, United Kingdom) (SIGCOMM '15). Association for Computing Machinery, New York, NY, USA, 537–550.
- [26] Radhika Mittal, Alexander Shpiner, Aurojit Panda, Eitan Zahavi, Arvind Krishnamurthy, Sylvia Ratnasamy, and Scott Shenker. 2018. Revisiting network support for RDMA. In *Proceedings of the 2018 Conference of the ACM Special Interest Group on Data Communication* (Budapest, Hungary) (SIGCOMM '18). Association for Computing Machinery, New York, NY, USA, 313–326.
- [27] ns-3 Project. 2025. ns-3: A Discrete-Event Network Simulator. <https://www.nsnam.org/>. Accessed: 2025-05-02.
- [28] NVIDIA. 2024. ZTR-RTT Congestion Control Algorithm Overview v1.0. <https://docs.nvidia.com/networking/display/ztrrttcongestioncontrolalgorithmoverviewv10> Accessed: 2025-05-16.
- [29] NVIDIA. [n.d.]. ConnectX-5 SmartNIC Adapter. <https://www.nvidia.com/en-gb/networking/ethernet/connectx-5/>. Accessed: 2025-05-28.
- [30] NVIDIA. [n.d.]. ConnectX-6 SmartNIC Adapter. <https://www.nvidia.com/en-gb/networking/ethernet/connectx-6/>. Accessed: 2025-05-28.
- [31] DPDK Project. 2025. Data Plane Development Kit (DPDK). <https://www.dpdk.org/>. Accessed: 2025-05-02.
- [32] SONiC Project. 2025. Software for Open Networking in the Cloud. <https://sonicfoundation.dev/>. Accessed: 2025-06-28.
- [33] Kun Qian, Yongqing Xi, Jiamin Cao, Jiaqi Gao, Yichi Xu, Yu Guan, Binzhang Fu, Xuemei Shi, Fangbo Zhu, Rui Miao, Chao Wang, Peng Wang, Pengcheng Zhang, Xianlong Zeng, Eddie Ruan, Ziping Yao, Ennan Zhai, and Dennis Cai. 2024. Alibaba HPN: A Data Center Network for Large Language Model Training. In *Proceedings of the ACM SIGCOMM 2024 Conference* (Sydney, NSW, Australia) (ACM SIGCOMM '24). Association for Computing Machinery, New York, NY, USA, 691–706.
- [34] Mubashir Adnan Qureshi, Yuchung Cheng, Qianwen Yin, Qiaobin Fu, Gautam Kumar, Masoud Moshref, Junhua Yan, Van Jacobson, David Wetherall, and Abdul Kabbani. 2022. PLB: congestion signals are simple and effective for network load balancing. In *Proceedings of the Conference of the ACM Special Interest Group on*

- Data Communication* (Amsterdam, Netherlands) (*ACM SIGCOMM '22*). Association for Computing Machinery, New York, NY, USA, 207–218.
- [35] Abhiram Ravi, Nandita Dukkipati, Naoshad Mehta, and Jai Kumar. 2024. *Congestion Signaling (CSIG)*. Internet-Draft draft-ravi-ippm-csig-01. Internet Engineering Task Force. <https://datatracker.ietf.org/doc/draft-ravi-ippm-csig/01/> Work in Progress.
- [36] Arjun Roy, Hongyi Zeng, Jasmeet Bagga, George Porter, and Alex C. Snoeren. 2015. Inside the Social Network's (Datacenter) Network. *SIGCOMM Comput. Commun. Rev.* 45, 4 (Aug. 2015), 123–137.
- [37] Cha Hwan Song, Xin Zhe Khooi, Raj Joshi, Inho Choi, Jialin Li, and Mun Choon Chan. 2023. Network Load Balancing with In-network Reordering Support for RDMA. In *Proceedings of the Conference of the ACM Special Interest Group on Data Communication* (New York, NY, USA) (*ACM SIGCOMM '23*). Association for Computing Machinery, New York, NY, USA, 816–831.
- [38] Erico Vanini, Rong Pan, Mohammad Alizadeh, Parvin Taheri, and Tom Edsall. 2017. Let It Flow: Resilient Asymmetric Load Balancing with Flowlet Switching. In *14th USENIX Symposium on Networked Systems Design and Implementation (NSDI 17)*. USENIX Association, Boston, MA, 407–420.
- [39] Zirui Wan, Jiao Zhang, Mingxuan Yu, Junwei Liu, Jun Yao, Xinghua Zhao, and Tao Huang. 2024. BiCC: Bilateral Congestion Control in Cross-datacenter RDMA Networks. In *IEEE INFOCOM 2024 - IEEE Conference on Computer Communications*. 1381–1390. <https://doi.org/10.1109/INFOCOM52122.2024.10621412>
- [40] Zilong Wang, Layong Luo, Qingsong Ning, Chaoliang Zeng, Wenxue Li, Xinchun Wan, Peng Xie, Tao Feng, Ke Cheng, Xiongfei Geng, Tianhao Wang, Weicheng Ling, Kejia Huo, Pingbo An, Kui Ji, Shidong Zhang, Bin Xu, Ruiqing Feng, Tao Ding, Kai Chen, and Chuanxiong Guo. 2023. SRNIC: A Scalable Architecture for RDMA NICs. In *20th USENIX Symposium on Networked Systems Design and Implementation (NSDI 23)*. USENIX Association, Boston, MA, 1–14.
- [41] William Won, Taekyung Heo, Saeed Rashidi, Srinivas Sridharan, Sudarshan Srinivasan, and Tushar Krishna. 2023. ASTRA-sim2.0: Modeling Hierarchical Networks and Disaggregated Systems for Large-model Training at Scale. In *2023 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*. IEEE, New York, NY, USA, 283–294.
- [42] Yunhong Xu, Keqiang He, Rui Wang, Minlan Yu, Nick Duffield, Hassan Wassel, Shidong Zhang, Leon Poutievski, Junlan Zhou, and Amin Vahdat. 2022. Hashing Design in Modern Networks: Challenges and Mitigation Techniques. In *2022 USENIX Annual Technical Conference (USENIX ATC 22)*. USENIX Association, Carlsbad, CA, 805–818.
- [43] Hong Zhang, Junxue Zhang, Wei Bai, Kai Chen, and Mosharaf Chowdhury. 2017. Resilient Datacenter Load Balancing in the Wild. In *Proceedings of the Conference of the ACM Special Interest Group on Data Communication* (Los Angeles, CA, USA) (*SIGCOMM '17*). Association for Computing Machinery, New York, NY, USA, 253–266. <https://doi.org/10.1145/3098822.3098841>
- [44] Jiao Zhang, Yuqing Wang, Xiaolong Zhong, Mingxuan Yu, Haoyu Pan, Yali Zhang, Zixuan Guan, Biyao Che, Zirui Wan, Tian Pan, and Tao Huang. 2024. PACC: A Proactive CNP Generation Scheme for Datacenter Networks. *IEEE/ACM Trans. Netw.* 32, 3 (Feb. 2024), 2586–2599. <https://doi.org/10.1109/TNET.2024.3361771>
- [45] Yanli Zhao, Andrew Gu, Rohan Varma, Liang Luo, Chien-Chin Huang, Min Xu, Less Wright, Hamid Shojanazeri, Mylène Ott, Sam Shleifer, Alban Desmaison, Can Balioglu, Pritam Damania, Bernard Nguyen, Geeta Chauhan, Yuchen Hao, Ajit Mathews, and Shen Li. 2023. PyTorch FSDP: Experiences on Scaling Fully Sharded Data Parallel. *arXiv:2304.11277 [cs.DC]* <https://arxiv.org/abs/2304.11277>
- [46] Yibo Zhu, Hagga Eran, Daniel Firestone, Chuanxiong Guo, Marina Lipshteyn, Yehonatan Liron, Jitendra Padhye, Shachar Raindel, Mohamad Haj Yahia, and Ming Zhang. 2015. Congestion Control for Large-Scale RDMA Deployments. In *Proceedings of the 2015 ACM Conference on Special Interest Group on Data Communication* (London, United Kingdom) (*SIGCOMM '15*). Association for Computing Machinery, New York, NY, USA, 523–536.

A HOPPER'S WORKFLOW

The workflow of Hopper is illustrated in Figure 7. Below, we explain the Hopper's workflow in detail.

- The congestion detection module actively monitors the RTT of the current path (P_3) (see Fig. 7(a)).
- When the RTT of the current path reaches the threshold th_prob , set to 12 μs in this example, the path probing module begins probing two alternative paths (P_1 and P_4) by creating two sets of QPs and binding them to different source ports (see Fig. 7(b)).
- If the RTT of the path exceeds the threshold th_cong , set to 14 μs in this example, the path switching module compares the congestion on the current path (P_3) with that of the probed alternatives (P_1 and P_4). If a considerably better alternative path is available, the flow is switched to it after a cautious delay proportional to the delay difference between the current path (P_3) and the alternative path (P_1), to minimize OOO packets at the receiver NIC (see Fig. 7(c)).
- Finally, the flow is switched to the selected alternative path (P_1) (see Fig. 7(d)).

B ALICLOUD STORAGE WORKLOAD

We present the results of ns-3 simulation for AliCloud [21] in Fig. 8 with the same parameters as for Meta Hadoop [36]. Under 50% load, our method provides up to 6% improvement for smaller flows and up to 19% for large flows on average, compared to CONGA and FlowBender. While it performs up to 6% worse than ConWeave for small flows, it still achieves up to 8% better 99th percentile performance than FlowBender and up to 16% better than CONGA. Under 80% load, our approach continues to outperform FlowBender by approximately 10% in average, while showing similar behaviour for the tail latencies.

C META HADOOP WORKLOAD

For completeness, we include the full flow size distribution in Fig. 9, which covers all flow sizes including those between 2KB and 49KB which were omitted in the section 4.1.2 for clarity.

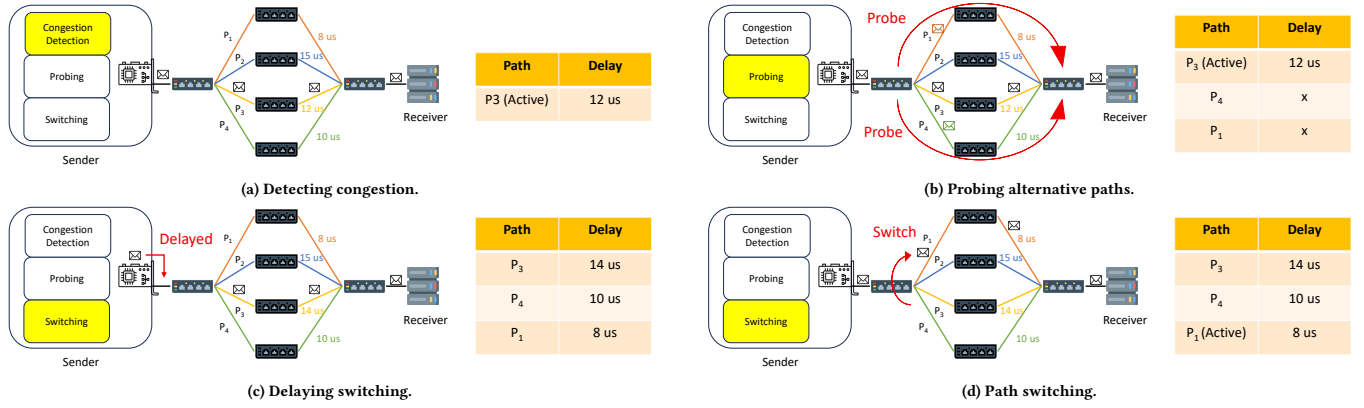


Figure 7: Hoppers workflow.

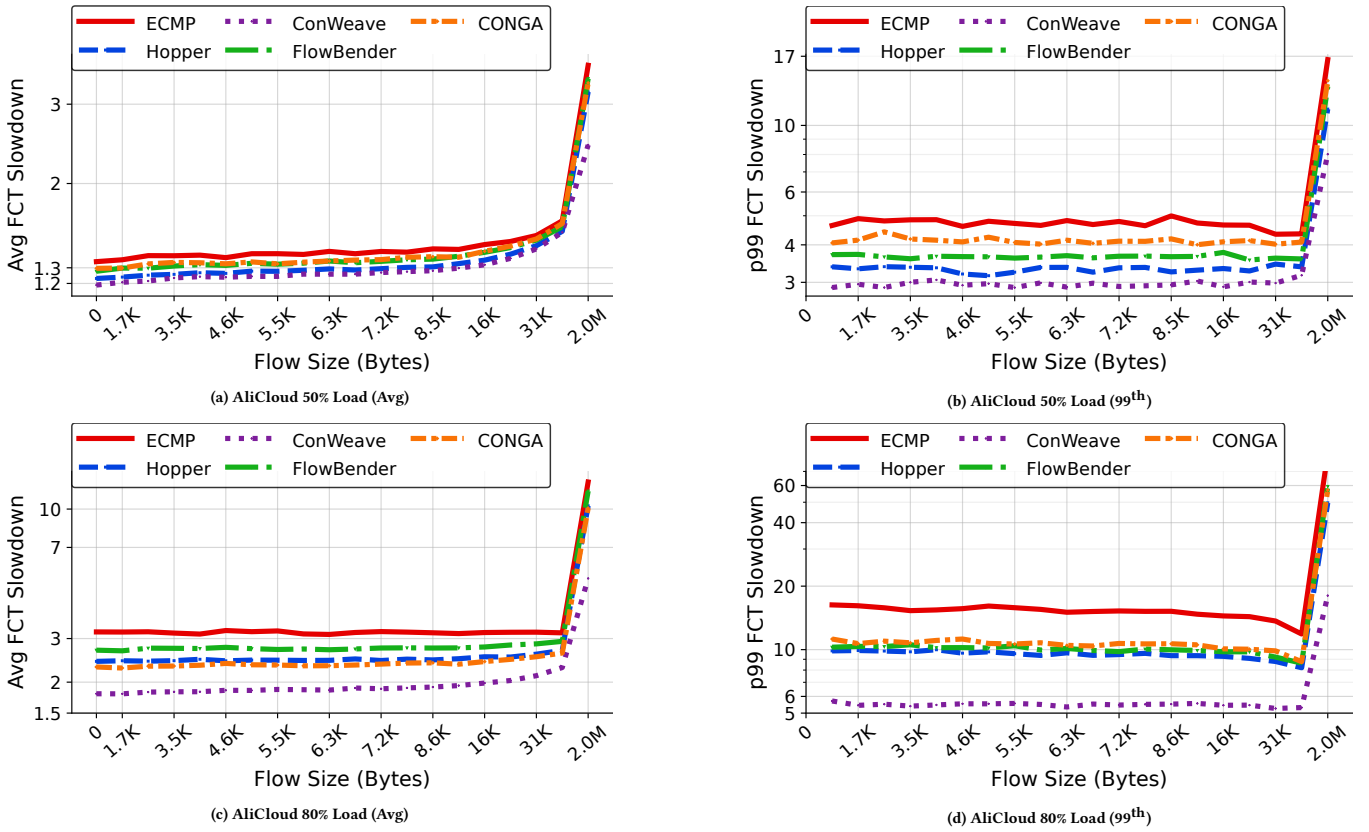
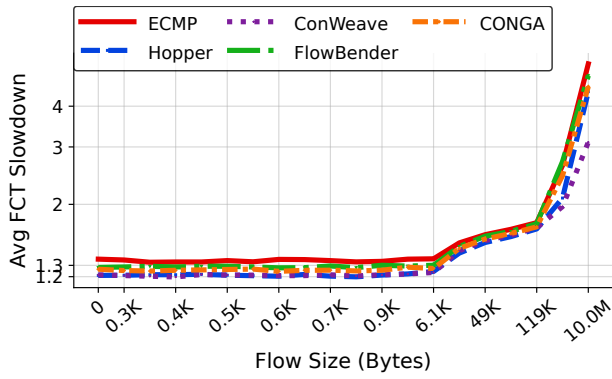
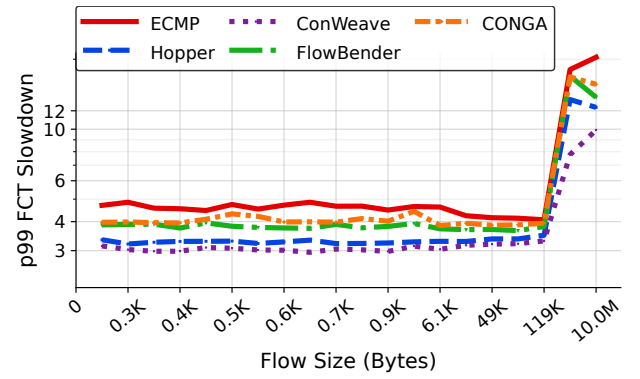
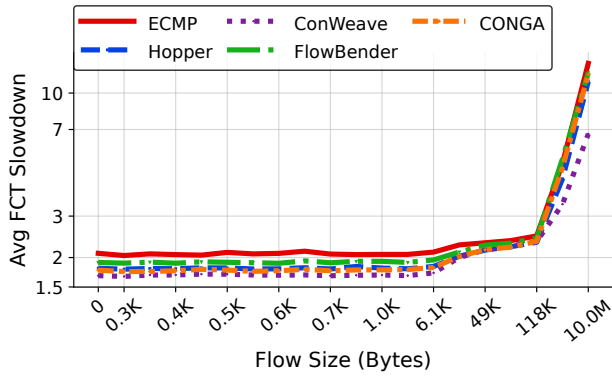


Figure 8: Average and tail FCT slowdown for AliCloud workload at 50% and 80% average network load.



(a) Hadoop 50% Load (Avg)

(b) Hadoop 50% Load (99th)

(c) Hadoop 80% Load (Avg)

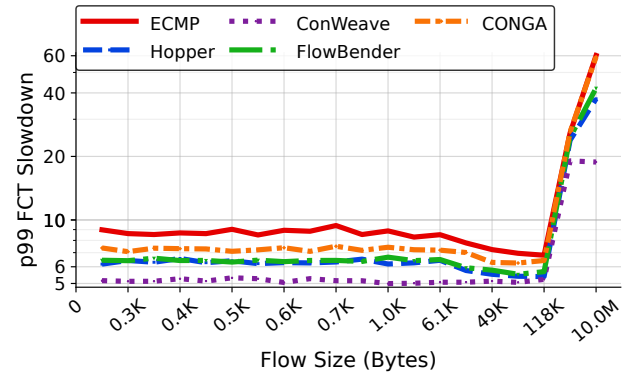
(d) Hadoop 80% Load (99th)

Figure 9: Average and tail FCT slowdown for Hadoop workload at 50% and 80% average network load.