# Impact of RoCE Congestion Control Policies on Distributed Training of DNNs

Tarannum Khan*, Saeed Rashidi†, Srinivas Sridharan‡, Pallavi Shurpali‡, Aditya Akella*, and Tushar Krishna†

*The University of Texas at Austin, Austin, USA

†Georgia Institute of Technology, Atlanta, USA

‡Meta, Menlo Park, USA

tarannum.khan@utexas.edu, saeed.rashidi@gatech.edu, ssrinivas@fb.com, tushar@ece.gatech.edu

*Abstract*—RDMA over Converged Ethernet (RoCE) has gained significant attraction for datacenter networks due to its compatibility with conventional Ethernet-based fabric. However, the RDMA protocol is efficient only on (nearly) lossless networks, emphasizing the vital role of congestion control on RoCE networks. Unfortunately, the native RoCE congestion control scheme, based on Priority Flow Control (PFC), suffers from many drawbacks such as unfairness, head-of-line-blocking, and deadlock. Therefore, in recent years many schemes have been proposed to provide additional congestion control for RoCE networks to minimize PFC drawbacks. However, these schemes are proposed for general datacenter environments. In contrast to the general datacenters that are built using commodity hardware and run general-purpose workloads, high-performance distributed training platforms deploy high-end accelerators and network components and exclusively run training workloads using collectives (All-Reduce, All-To-All) communication libraries for communication. Furthermore, these platforms usually have a private network, separating their communication traffic from the rest of the datacenter traffic. Scalable topology-aware collective algorithms are inherently designed to avoid incast patterns and balance traffic optimally. These distinct features necessitate revisiting previously proposed congestion control schemes for general-purpose datacenter environments. In this paper, we thoroughly analyze some of the state-of-the-art RoCE congestion control schemes (DCQCN, DCTCP, TIMELY, and HPCC) vs. PFC when running on distributed training platforms. Our results indicate that previously proposed RoCE congestion control schemes have little impact on the end-to-end performance of training workloads, motivating the necessity of designing an optimized, yet low-overhead, congestion control scheme based on the characteristics of distributed training platforms and workloads.

*Keywords*-distributed training; collective communication; network congestion control; RDMA over Converged Ethernet (RoCE)

## I. INTRODUCTION

Deep Neural Network (DNN) models are gaining significant attention due to their wide applicability across different domains such as vision [1], [2], [3], language modeling [4], sensor networks [5], and recommendation systems [6]. DNN workloads first need to be trained on existing samples to achieve high accuracy. However, as the number of training samples and DNN model sizes increase, training is becoming increasingly challenging, requiring several days/weeks to be trained [7], [8], [9]. This is especially true for DNNs such as recommendation models [6] that need to be constantly trained using new data samples produced daily.

Distributed training is used in practice today as a solution to reduce the training time by distributing the training task across multiple accelerators, aka *Neural Processing Units* (NPU - a term that abstracts any specific accelerator type (e.g., TPU, GPU, FPGA)). However, distributed training comes at the expense of communication overhead between NPUs to synchronize model gradients and/or activation depending on the *parallelization strategy* of the distributed training. This communication is usually handled through a series of synchronous and multi-peer patterns, called *collective communication patterns* (e.g., All-Reduce) [7], [10]. Scalable topology-aware collective algorithms are designed to avoid in-cast patterns and balance traffic optimally.

To achieve maximum performance, in recent years specialized distributed training platforms have been built, both by cloud providers [11], [12], [13] and accelerator vendors [14], [15]. These platforms have distinct characteristics that set them apart from the conventional data-center environments:

First, compared to the traditional data-center setting that relies on commodity hardware, distributed training platforms are built using high-end compute and network components [12], [16], [17]. In such an environment, it is typical for each NPU to have a dedicated high-BW Network Interface Card (NIC), supported by direct remote memory access (RDMA) to distant NPUs. This provides enormous network BW per server node, given that each server node hosts multiple NPUs (up to 16 per server node [16]).

Second, to mitigate the communication overhead, distributed training platforms employ dedicated networks that separate training traffic from the rest of the datacenter traffic [12]. In this case, each server node has a dedicated NIC per CPU socket to connect to the datacenter network, while NPUs across different server nodes form their private network using their dedicated NICs [12].

Third, due to the growing size of DNN models and the significantly growing size of training datasets that are generated daily, training platforms are often scheduled to perform only one training job at a time for the critical DNN workloads (e.g., recommendation models) to minimize the training time and enable training over larger datasets [18].

As a result of these unique characteristics, it is necessary

to revisit the networking stack and identify whether current state-of-the-art networking protocols are optimal for such platforms. In particular, we focus on *RDMA over Converged Ethernet* (RoCE) protocol [19] due to its compatibility with current Ethernet-based fabric and widespread usage on distributed training platforms [12]. Recent works have shown the importance of congestion control on RoCE to achieve maximum performance [20], [21], [22], [23]. This is because the RDMA protocol is more efficient on lossless networks, which is not natively supported on Ethernet-based fabrics [19]. To address this issue and achieve near-lossless network guarantees, baseline RoCE enforces congestion control at the link layer through the Priority Flow Control (PFC) mechanism. In this case, once the receiver-side buffer occupancy crosses a threshold, the receiving NIC generates a PAUSE frame and sends it to the sender who stops sending packets until further notified by the receiver.

However, as previous works have shown, the PFC mechanism suffers from many drawbacks in conventional data-center environments, including unfairness, head-of-line-blocking, and deadlock [20]. Therefore, various works have proposed additional NIC- and network-based congestion control techniques on top of PFC enforced to achieve high-performance and lossless features while minimizing the PAUSE frames generated by PFC. [20], [21], [22], [20], [23].

In this paper, we study the effect of these state-of-the-art congestion control mechanisms – proposed for general data-center workload platforms/environments – and compare them with baseline PFC in the context of distributed training workloads and platforms. We study the impact on end-to-end training time and explain the observations.

We make the following contributions:

- To the best of our knowledge, this is the first work that evaluates the effect of different congestion control schemes on distributed training.
- We developed a simulator using ASTRA-Sim [24] and NS3 [25]. We model state-of-the-art training workloads and platforms in ASTRA-Sim [24] and integrate with NS3 [25] to model the RoCE network and its different congestion control schemes.
- We provide a detailed analysis of the effect of each state-of-the-art congestion control scheme (i.e., Baseline PFC, DCQCN, DCTCP, Timely, and HPCC) for both single collective communication micro-benchmarks and end-to-end training time of the DLRM workload [6].
- We show that different state-of-the-art RoCE congestion control schemes have little impact on the end-to-end training performance.
- Based on our analysis, we provide directions for designing an optimized yet low-overhead congestion control scheme tuned for distributed training.

## II. BACKGROUND

### A. Collective Communication Patterns

Collective communications are the most frequent communication pattern observed in distributed training [7], [9], [26]. Collective communication patterns refer to a set of collaborative and synchronous data exchange between NPUs for a specific purpose [18], [7]. The most frequently used pattern in distributed training is *All-Reduce*. In All-Reduce, each NPU initially has a data buffer, and the end goal is to *reduce* the data buffer of all NPUs and replicate the results on all NPUs [27], [28]. *All-Reduce* can be broken into executing two consecutive communication steps: i) *Reduce-Scatter* followed by ii) *All-Gather* [9]. In the Reduce-Scatter phase, each NPU obtains a portion of the globally reduced data buffer, while in All-Gather, each NPU broadcasts their data buffer to all other NPUs [18].

In addition to All-Reduce, All-To-All is another collective pattern observed for some training workloads such as DLRM recommendation models [6]. In All-To-All, each NPU sends a specific portion of each data buffer to each NPU.

### B. Collective Communication Algorithms

**All-Reduce.** To handle the collective patterns discussed in Section II-A, many different algorithms have been proposed in recent years. For example, basic All-Reduce algorithms include: ring-based [27], tree-based [28], halving-doubling [29], etc. Each basic algorithm is optimal for certain physical topologies. For example, the ring-based algorithm is optimized for physical ring topology, while direct All-Reduce works well for physical switch-based topologies. [18].

In the basic direct All-Reduce algorithm among P NPUs, the data buffer is split into P segments, and each NPU sends its i'th segment ($0 \leq i < P$) to the NPU id #i at the same time. Therefore NPU id #i receives all i'th segments of all NPUs, which then reduce with its local i'th segment (Reduce-Scatter). In the next phase, each NPU broadcasts its locally reduced segment to all other NPUs at the same time (All-Gather) [18].

For heterogeneous network topologies with multiple levels of the network, multi-stage hierarchical All-Reduce algorithms have been proposed [9], [28]. In this case, All-Reduce is broken into a sequence of Reduce-scatter operations starting from the first level and going all the way to the last network level. This is then followed by issuing All-Gathers in the reverse order of Reduce-Scatter. The Reduce-Scatter/All-gather algorithm for each stage is a basic collective algorithm discussed earlier (e.g., ring-based, halving-doubling, etc.). Doing so reduces network traffic as data goes to the next network level, which is desired since later network levels have less bandwidth than the first levels.

**All-To-All.** Like All-Reduce, the All-To-All pattern can also be implemented in multiple ways. However, the most common algorithm is direct All-To-All, where each NPU

sends its messages to other destination NPUs at the same time [18], [30], [31]. The reason is that, unlike All-Reduce, hierarchical All-To-All does not reduce the network traffic as data goes to the next network level.

### C. DLRM Workload

Recommendation models are one of the most important class of DNN workloads that need to be trained constantly with a huge amount of user data generated daily. Therefore, they serve as a suitable target for our studies since they require maximum training performance to keep up with the growing user training data [12].

DLRM [6] is the most common DNN model used for recommendation systems. DLRM combines user-obtained sparse features (e.g., pages the user previously liked) and dense features (e.g., age) and produces the *Clickthrough rate (CTR)* for a user to view a specific Ad. It consists of three main parts: bottom MLP, embedding layer, and top MLP. Bottom MLP is used to process dense features while embedding tables represent sparse features. Sparse features are looked up in the tables, and then the output of embedding tables interacts with each other (e.g., using dot product) to generate the output of the embedding layer. The embedding layer and bottom MLP layers can run in parallel, and their outputs are then concatenated to feed the top MLP layers to predict CTR.

The two most common parallelization strategies in DNN training is *data-parallel* and *model-parallel*. In data-parallel, the model is replicated across NPUs, and each NPU works on separate mini-batches. This strategy requires an All-Reduce between model gradients computed by all NPUs before updating the model and starting the new training iteration. In model-parallel, the model is split across NPUs, but all NPUs work on the same minibatch. Therefore it requires collective communications between NPUs to exchange output activations/input gradients of each model-parallel layer during forward-pass/backpropagation. The exact model-parallel communication type depends on the layer type.

The standard parallelization strategy of DLRM is to distribute the MLP layers in a data-parallel manner, while the embedding tables are split across all NPUs due to their huge sizes, forming a model-parallel split [6], [12]. Therefore, DLRM requires All-Reduce collectives for data-parallel MLP layers, and All-To-All for forward-pass/backpropagation of embedding layer [6].

### D. RoCE Congestion Control Schemes

We evaluated our target workloads on the following congestion control policies:

*1) PFC only:* In this case, we did not use any congestion control algorithm and used only pause frames to control the congestion [19]. When the bytes in the queues are more than a threshold value, pause frames (PFCs) are triggered,

and when these PFCs reach the sender side, the sender stops sending more packets until the sender receives resume frames sent by the switch.

*2) DCQCN:* DCQCN [20] is inspired by QCN [32] and DCTCP [33]. Once the queue length surpasses the ECN threshold, a congestion notification packet is generated from the receiver of the ECN-marked packet to inform the sender to reduce the flow rate. If the sender receives no feedback for a fixed time unit, the sender starts to increase the rate. Initially, it increases the rate rapidly in the fast recovery phase, and then it additively increases the rate. DCQCN starts at line rate. DCQCN has many parameters that need to be tuned for better performance.

*3) DCTCP:* DCTCP [33] achieves low latency and high throughput requirements by reacting to congestion in proportion to congestion. DCTCP uses a simple marking scheme at switches. It marks packets by setting the ECN flag at the switches if the queue occupancy exceeds the threshold. The sending window is reduced at the sender side depending on the reduction factor, alpha. DCTCP was originally proposed for the TCP network. We used the same algorithm developed for datacenter TCP traffic but used it for RoCE v2 as is done in HPCC [23] with the starting rate as line rate so that it can be equivalently compared with other congestion algorithms which are starting at line rate.

*4) TIMELY:* TIMELY [22] is a congestion control mechanism that measures round trip time (RTT) delay and accordingly decides to reduce or increase the rate at the sender side. If the RTT is less than Tlow, the rate is increased additively. If the RTT is greater than Thigh, the rate decreases multiplicatively, and if RTT is between Tlow and Thigh, the rate increases or decreases based on the gradient.

*5) High Precision Congestion Control (HPCC):* Earlier algorithms were doing congestion control based on ECN or delay. But, HPCC does congestion control through the in-network telemetry, INT feature of its switching ASIC where it inserts information such as timestamp, queue length, etc. When the receiver receives its packet, it copies this metadata to the ack and sends the ack packet to the sender to take action depending on the metadata information in the ack packets. HPCC also controls inflight bytes by using a window and varying it according to the INT header information. HPCC does some more interesting optimization to monitor congestion closely [23].

*6) HPCC-PINT:* In HPCC, there is an INT overhead for every packet as each switch adds this information to the packet. A data center topology with five hops with the only addition of two values per switch requires 48 bytes of overhead or 4.8 percent of a 1000 bytes packet. The HPCC problem of overhead per packet is solved using PINT with HPCC, which uses only 8 bits. HPCC-PINT [34] achieves this by not giving feedback per packet but depending on a parameter. Hence feedback can sometimes be delayed, especially for shorter flows.
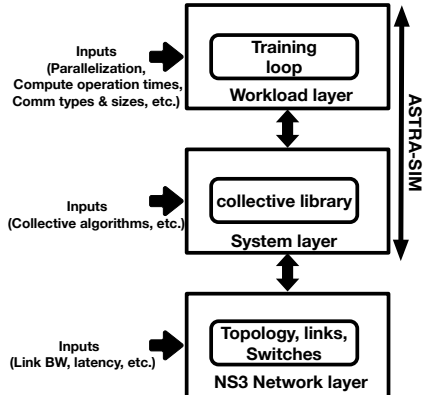
Figure 1: The Simulation methodology architecture

## III. METHODOLOGY

### A. Simulation Methodology

We use the ASTRA-SIM simulator [35][24] to model the communication scheduling of deep learning (DL) training workloads (both individual collectives and actual DL models). ASTRA-SIM is an open-source cycle-level simulator for DL training platforms developed by Georgia Tech, Meta, and Intel [24], validated against current training platforms. Fig. 1 shows an overview of ASTRA-SIM. It comprises three layers. The *workload layer* takes the DNN model descriptions as an input, in terms of compute times and communication operations per layer, and simulates the training loop. In this work, we use NVIDIA V100 GPU profiling to obtain compute times. For simulating the communication operations, the workload layer calls the *system layer* communication APIs that implement diverse collective communication patterns via various algorithms (e.g., ring-based, halving-doubling, and so on). To accurately model the collective communications, the system layer breaks the algorithm into a series of send/recv operations that should be modeled via the *network layer*. ASTRA-sim provides a network API [18] to plug in diverse network backends. We integrated the Alibaba NS3 simulator [25] as the network layer in this work. NS3 models different state-of-the-art RoCE congestion control (CC) schemes including: Baseline PFC [19], DCQCN [20], Timely [22], DCTCP [33] and HPCC [23]. We use these CC algorithms for our comparative analysis. To the best of our knowledge, this is the first integrated simulator that models distributed training workloads on RoCE networks.

### B. Target Platforms

We model the distributed training platform similar to Mudigere et al. [12]. Fig. 2 presents the topology architecture. In this case, each server nodes host 8 GPUs, locally connected to six NVSwitches (scale-up network) within a server node. Each GPU has a dedicated NIC connecting it to the top-of-the-rack (TOR) switch (scale-out network). There are two such server nodes within a rack. The inter-rack connectivity is enabled by connecting TOR switches to spine switches as shown in Fig. 2. The connectivity is a 1:1 full subscription
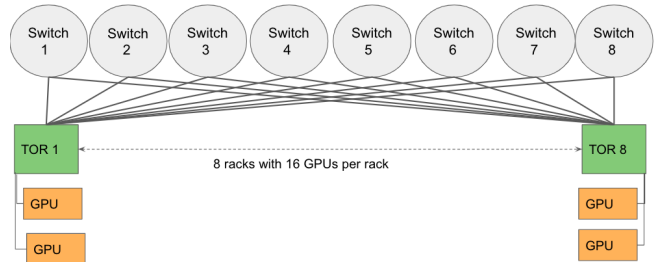


Figure 2: The two-level Clos topology commonly used in datacenters.

Table I: Topology design parameters

| Parameter | Value |
|---|---|
| Scale-up (NVlink) link BW | 200 GBps (total) |
| Scale-up (NVlink) link latency | 25 ns |
| NIC-to-TOR link BW | 200 Gbps (total) |
| NIC-to-TOR link latency | 500 ns |
| TOR-to-Spine link BW | 200 Gbps (total) |
| TOR-to-Spine link latency | 500 ns |
| TOR Switch buffer size | 32 MB |
| Spine switch buffer size | 32 MB |

non-blocking CLOS topology. This paper simulates up to sixteen racks with eight spine switches.

Table I summarizes the topology design parameters. Note that in this design, the GPU network is separated from the general datacenter traffic for maximum training performance[1][12].

### C. Target Workloads

We use incast, single collective micro-benchmark workloads (i.e., All-Reduce and All-To-All), and real end-to-end training workloads to compare different CC algorithms. The goal is to first understand CC behaviours for simpler communication patterns and then extend it to the full end-to-end training performance. For real training workloads, we pick distributed training of the DLRM [6] recommendation model since it is one of the critical workloads that need constant training over newly generated samples and hence, requires maximum training performance [12]. The DLRM model size is similar to [18]. Table II shows our DLRM description.

**Multi Tenancy.** As discussed in Section I, to model the critical training workloads, we assume only a single training workload is under execution on our target platforms. This is inline with many real system scenarios where maximum performance is required to keep pace with the constant need to train the network with newly generated training data produced daily [12], [18], [7], [36].

### D. Target Collective Communication Algorithms

For All-Reduce, we use two different versions of All-Reduce: i) *1D*: which is a basic direct All-Reduce algorithm,

---

[1]In this case, the connectivity of each server node to the rest of the datacenter is enabled via the CPU NIC (not modeled in this paper).

Table II: DLRM Model parameters

| Model Parameters | Value |
|---|---|
| Size of embedding data-type | 16 bits |
| Pooling factor | 60 |
| Top MLP layers | 10+2 |
| Bottom MLP layers | 5+2 |
| Dense features | 1600 |
| Top MLP layer size | 2048 |
| Bottom MLP layer size | 1024 |
| Sparse features | 64 |
| Embedding dimension | 64 |

and ii) *2D*: which is a BW-efficient hierarchical algorithm as described in Section II. This means that on our target topology, an All-Reduce is broken into: Reduce-Scatter within each server node through scale-up network (NVLink), followed by Reduce-Scatter between GPUs of the same id across different server nodes through scale-out NICs, followed by All-Gather in the reverse order (four steps in total). Since our topology is switch-based at all levels, the basic All-Reduce algorithm on each stage is a direct algorithm as it is shown to be efficient on switch-based networks [18].

For All-To-All, we perform the direct All-To-All algorithm described by Rashidi et al. [18] and Section II-A.

Also, to utilize all network levels simultaneously, we break each collective into four equally sized chunks and then process them in a pipeline manner, similar to [37].

### E. Metric of Evaluation

We use various metrics to evaluate the performance and efficiency of CC algorithms for microbenchmark and real workload cases. We mainly focus on network stats such as end-to-end completion time, switch buffer occupancy, and #of PFC PAUSE frames for microbenchmark communication workloads.

For real workloads, we focus on mainly workload-related stats such as end-to-end training iteration time that can be decomposed into *total compute times + exposed communication time*. Exposed communication time is when the training loop is forced to stop because it needs to wait for the communication to be finished. In other words, it shows the amount of communication time that cannot be overlapped with compute.

### IV. RESULTS & ANALYSIS

#### A. Single-switch Incast micro-benchmark

For studying incast, we used a simple topology where 8 GPUs are connected to a switch. The switch has a buffer size of 32 MB. All the GPUs are connected to a switch with a link of bandwidth 200 Gbps and 500 ns latency. Seven of the GPUs are sending 10MB data to GPU 0. In all the figures in result, queue length in bytes is the sum of bytes in all the queues in a switch at a time. We study how the egress buffer queue is varying in all the algorithms. Note that in this case, the buffer size will be saturated first when compared with the link connecting GPUs and the switch as queue buffer is

a scarce resource compared to the link bandwidth. Hence utilization of the queue such that we are not underutilizing the bandwidth becomes an important criteria for different congestion control algorithm to perform well. We observe similar performance for all CC algorithms.

*1) PFC only:* For the case where only PFC [19] is enabled without any additional congestion control, we see there is a queue build up throughout the flow time as shown in Fig. 3. This is because, once the queue is build and reaches the threshold value of a switch queue, PFCs are triggered and whenever the queue length is dropped below the threshold, resume is triggered which allow again more packets to pass through the link and into the switch. Hence, maximum queue length is utilized and also the link bandwidth is utilized efficiently. In case of using only PFC, there are no complex computation which needs to be done like in other congestion control policies. Hence, this also saves some time. But this advantage, comes with a cost of a lot of PFC's to get triggered. There are a lot of problems with PFC like head-of-line blocking, deadlocks and unfairness [20]. Because of these reasons, various congestion control were introduced to avoid PFC production as much as possible.

*2) DCQCN:* The first congestion control algorithm, we will study is Datacenter QCN, DCQCN [20] which can be seen in Fig. 3. DCQCN is a fluid model with various parameters in switch to be set. We tuned the parameter in every experiment to test against the best case of DCQCN. In this case, there are no PFC generated and we get better link utilization with less queue buffer usage than the previous case where a lot of PFCs were generated. Hence, we have low latencies and also eliminating the problem which comes with PFC. In the incast case here, first time there is a shoot up in the switch queue, as there is a delay in the initial rate reduction as the sender will have to wait for the first congestion notification packet to be received from the switch. After, that we can see the queue buffer has stabilised to be within a range.

*3) DCTCP:* The next algorithm we studied is DCTCP [33]. As we can see in Fig. 3, the queue is being built up initially but then due to the ecn notification received from the switches the sender decreased there window size and the eventually the queue size becomes very small. Here, also the PFC trigger point is not reached and no PFCs are being produced.

*4) TIMELY:* As we can see in Fig. 3, TIMELY reduces the rate heavily initially. As TIMELY receives delayed RTT initially because of the switch queue buildup, it starts aggressively reducing the rate of the senders. As the rate of all the seven GPUs sending to the zeroth GPUs get reduced together, the switch buffer usage becomes very low. With the low queue buffer usage, the bandwidth is also underutilized and the latency is worst for TIMELY. We tried various parameters but we were getting almost the same latency, so we used the values provide in the TIMELY paper. Hence,
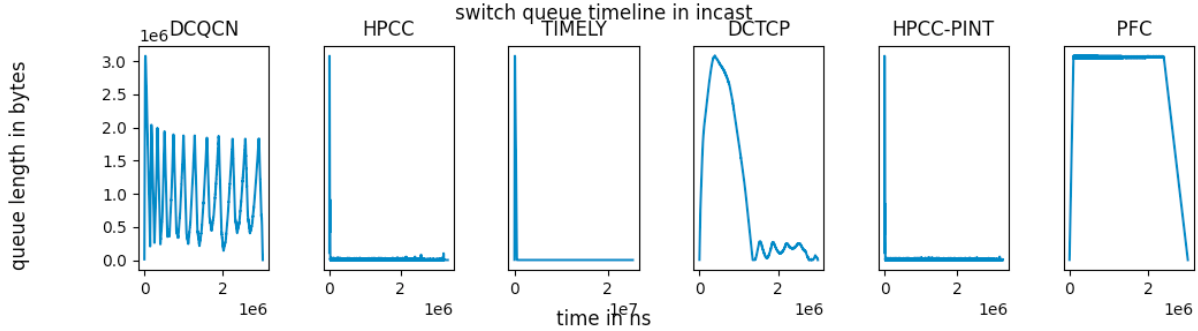
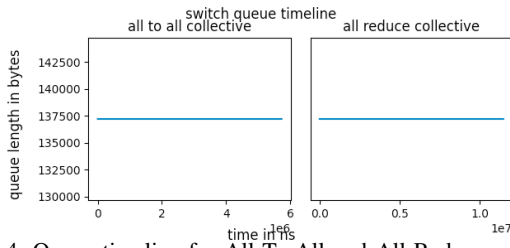Figure 3: Queue length timeline in incast



Figure 4: Queue timeline for All-To-All and All-Reduce collective



Figure 5: Effect of ECMP as observed from queue length timeline for different Spine switches for the same All-To-all Collective.

for large flow, very less queue buildup may not be a good option as there are no shorter flows here which will get affected by the queue build ups. As we can see in DCQCN, which has the least latency, a queue build up which is also getting consumed quickly is an ideal buffer pattern for large flows, so that the bandwidth can also be utilized in this scenario and PFCs are not produced. In TIMELY too for this case, PFCs are not generated.

*5) HPCC:* As we can see in Fig. 3, HPCC through the INT Header is continuously monitoring the queue length when it start receiving it acks and hence, continuously starts decreasing the window size until the queue length of the switch is minimal, which is the aim of HPCC. In this case, also no PFCs are produced. It is worth noting every packet there is an INT overhead as each switch adds this information into the packet. Hence, we are actually transferring more data than in other congestion control schemes and this may increase flow completion time. It also increases processing latency at switches [34].

*6) HPCC-PINT:* The solution to the HPCC problem of overhead per packet is solved using PINT along with HPCC which uses only 8 bits and have similar results and time series of queue length when compared with HPCC as seen in Fig. 3. HPCC-PINT achieves this by not giving feedback per packet but depending on a parameter, resulting in delayed feedback. This delayed feedback will worsen results for shorter flows but since in our study we are focusing on collectives which are generally longer flows, the results are better as the overhead per packet is reduced.

### B. Single-switch Collectives Micro-benchmark

We studied All-Reduce and All-To-All workload for one switch connected to 8 GPUs and scaled it to one switch
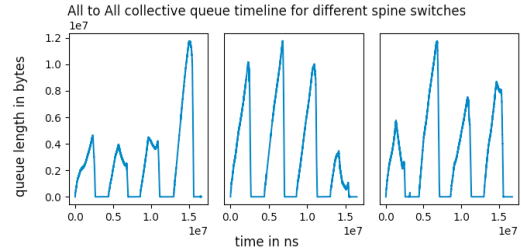
connected to 128 GPUs. We use the bandwidth of 200 Gbps and 500 ns latency and a switch of buffer size 32 MB. We also studied the impact of different congestion control algorithms for these workloads. Due to the nature of how these collectives are designed, we observed there was no congestion in this network and all the congestion control algorithms have similar completion times.

Let's consider All-To-All where each node is sending data to all other nodes. If the total collective data to be communicated is 10MB, then in the case where there are 8 GPUs, each GPU will send 1.25MB data to each of the other GPU. In the case of incast all of the 7 GPUs were sending data to the same queue of the switch. Thus, one switch queue was receiving 7*10MB of data. In All-To-All, as the data is distributed among different GPUs, it is also distributed among different switch queues. Hence, each of the switch queue will received 8 * 1.25MB of data from the 8 GPUs which is not enough to cause congestion, therefore all the algorithms produce similar latency. We have not added simulation results for less than 128 GPUs considering space limits as we see similar results as in Fig. 4, when sending 10 MB of data from all the nodes in All-Reduce and All-To-All for 8 GPUs. We observe that the sum of all the switch queue lengths is constant with respect to time showing no queue buildup and no congestion for all the congestion control policies. When this is scaled from 8 GPUs to 128 GPUs and the collective communication is scaled up to 128MB, similar results are obtained due to the same pattern of data distribution. The same pattern is observed in All-Reduce as well. Hence, in Fig. 4, we can see that there is no queue buildup and no congestion. As there were no congestion,
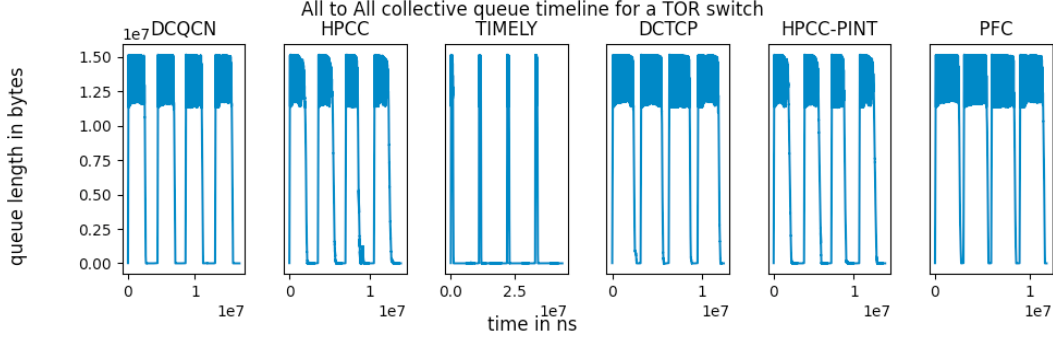
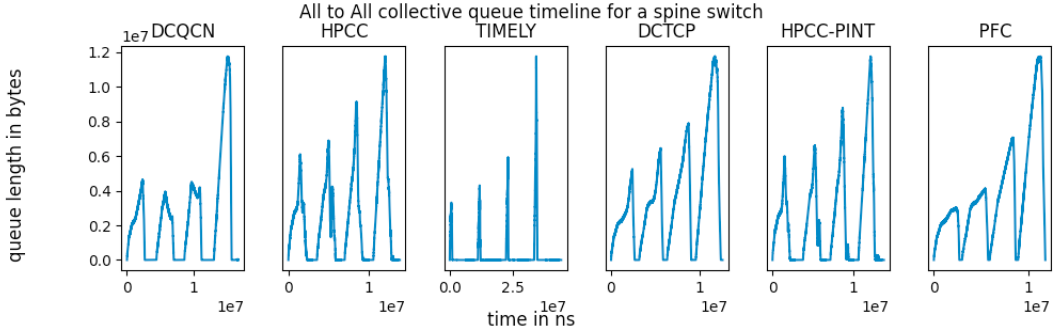Figure 6: Queue timeline for a TOR switch for different CC for All-To-All collective



Figure 7: Queue timeline for a Spine switch for different CC for All-To-All collective

we also observed there were no PFC generation for all the congestion control algorithms.

### C. Two-level CLOS topology

In this section, we study the effect of congestion on the common topology used in the datacenter, which is a two-level CLOS topology. This topology is modeled similar to [12] and shown in Fig. 2. We studied All-Reduce, All-To-All, and the DLRM workload for different congestion control policies in this topology setting.

*1) All-To-All Collective:* Fig. 6 and Fig. 7 depicts queue timeline for ToR and Spine switches for all the congestion control algorithms for All-To-All collective. We observe four peaks in Fig. 6 and Fig. 7 as we have split the data into four chunks which means we are completing four All-To-All flows from all the GPUs, one by one. As each chunk issues, all sends in one burst and then waits for completion of all its messages, therefore, we see four queue length peaks in Fig. 6 and Fig. 7. Fig. 5 depicts the queue timeline for three different Spine switches for the same All-To-All collective. From Fig. 5, we can observe that the Spine switches have different queue build-ups simultaneously for the same All-To-All collective flow. As the same Spine switch can be responsible for sending data to different TORs depending on the ECMP hash, in some of the switches, there could be more queue build-up compared to other switches at a point in time, as seen in Fig. 5. Hence, we can see queue build-up and congestion at different times in different Spine switches, which is depicted in Fig. 5 from the perspective of

different switches and Fig. 7 from the perspective of different congestion control algorithms. For TOR switches in those communication periods, the queue lengths are in a constant range as seen in Fig. 6 which indicates that there is not much congestion in the TOR switches.

All the congestion control algorithms have a similar queue build-up and latency except for TIMELY, which can be seen in Fig. 8 where we perform an All-To-All collective for 128 MB data. Similar results were observed for other data sizes as well. TIMELY being a rate-based algorithm, after receiving RTT delay way more than expected initially, it reduced the rate heavily and hence started underutilizing the bandwidth. We also suspect another reason for TIMELY poor performance could be non-optimal parameters. We used the parameters mentioned in the TIMELY paper and played with parameters close to the ones mentioned in the TIMELY paper. However, we still observed poor performance when using TIMELY compared to other algorithms.

Overall, while the PFC-only scheme generates the maximum number of PFCs, as shown in Fig. 9, we observe it achieves the equal or best performance compared to other CC schemes (except TIMELY). Because we have long flows, PFC does not degrade the BW utilization.

*2) All-Reduce Collective:* As in All-To-All, All-Reduce also has congestion because of ECMP, as explained in the previous section. For All-Reduce, we see a similar queue buildup for TOR and spine switches for all the congestion control algorithms as in All-To-All collective. For All-Reduce completion time as seen in Fig. 8, we ran 128 MB All-Reduce
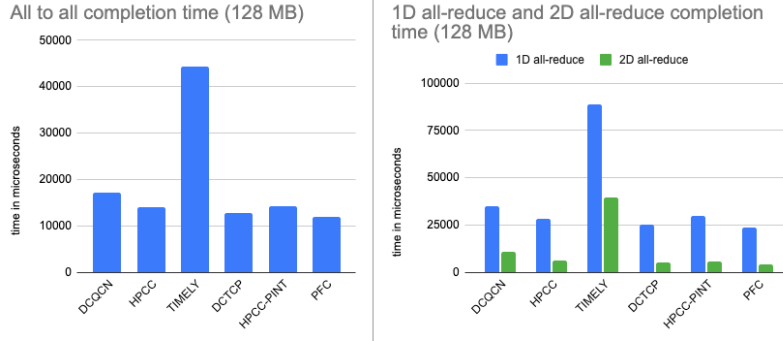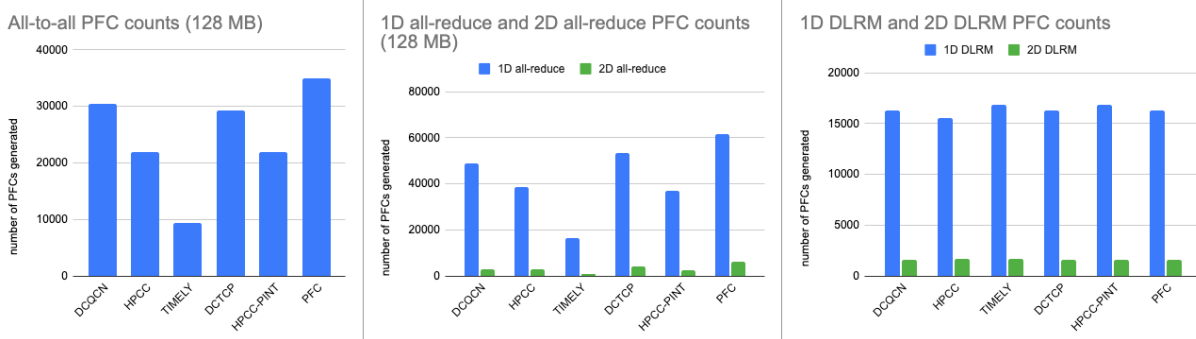
Figure 8: Completion time for collectives



Figure 9: PFC counts for workloads

for 128 GPUs. Similar results were observed for other data sizes as well. Similar to all-to-all, in all-reduce also because of the same reasons, the performance of TIMELY is the worst among all congestion control algorithms. All other congestion control algorithms have similar performance compared to only PFCs with no congestion control algorithm.

Additionally, we compared how 1D and 2D collectives fare in this topology set. 1D All-Reduce is doing All-Reduce throughout 128 GPUs. 2D All-Reduce means first doing All-Reduce in each server node and then doing All-Reduce among other GPUs in the servers, hence sending less data into the TOR and Spine switches in 2D All-Reduce compared to 1D All-Reduce. When comparing 1D collective versus 2D collective, we can see a significant reduction in completion time from Fig. 8 for all the congestion control policies. Also, from Fig. 9, we also observe a significant reduction in the count of PFCs in 2D collective versus 1D collectives as we are sending less data into the network.

### D. Real Workload: DLRM Results

In this section, we show the effect of different CC algorithms on the end-to-end training time of the DLRM model described in Section III-C. Fig. 10 shows one training iteration time, decomposed into total compute and total exposed communication, for different CC algorithms and two variations of the All-Reduce algorithm (i.e., 1D and 2D). In total, each training time of DLRM issues 109.5 MB All-Reduce and 8 MB All-To-All for MLP and embedding layers, respectively.
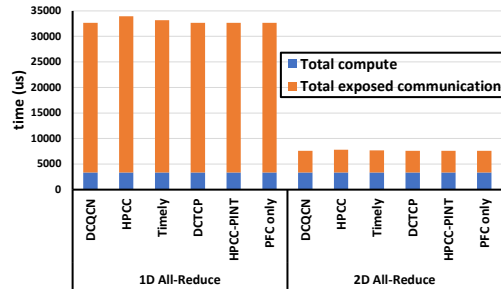


Figure 10: The end-to-end total compute and total exposed communication for 1 training iteration of DLRM model running on 128×V100 GPUs.

As can be seen, choosing the hierarchical 2D All-Reduce algorithm results in much better performance due to better utilizing local High-BW NVlinks and sending less data through NICs. However, current CC algorithms show a slight variance in end-to-end performance with each algorithm. More specifically, the end-to-end performance difference is less than 4% across all CCs, and No PFC gives the best performance results along with DCQCN, DCTCP, and HPCC-PINT. HPCC works worst in 1D and 2D All-Reduce cases mainly due to its added header overheads that reduce the *network goodput*. The performance of TIMELY is not as poor as in earlier experiments, as the same parameters are almost optimal in this scenario.

### E. Discussion

Scalable topology-aware collective algorithms are inherently designed to avoid in-cast patterns and balance traffic optimally. Our results show that baseline PFC is equal to other proposed congestion control (CCs) in terms of end-to-end performance. This means that state-of-the-art CC algorithms introduce considerable overhead at the endpoints, requiring frequent network feedback but offering little performance improvements (for the single workload scenario). Note that compared to the general datacenters, the overhead of CC algorithms are much higher in distributed training since they use high-performance network components (e.g., 400 Gbps NICs [38]) and also host multiple NICs per server node, wasting a lot of cycles at the endpoints to run CC algorithm. The only benefit of proposed CCs over the baseline PFC is that reducing the number of PAUSE frames minimizes the chance of PFC deadlocks that can rarely happen and halt the network. Additionally congestion control algorithms like DCQCN and TIMELY require extensive hyperparameter tuning and may not be optimal for all the flows. Additionally, tuning the congestion control hyperparameter before running every deep learning workload is not a feasible solution. Hence, we conclude that an optimized CC mechanism for distributed training can be much simpler than other complicated (with high overhead) CCs. The only feature optimized CC needs to provide is to prevent rare event deadlocks. This can be done by setting the congestion window to minimize buffer usage and hence minimize the PFC and chance of deadlocks. Fortunately, the communication patterns of distributed training are deterministic and repeated for each training iteration. Therefore an optimized CC can be a very low overhead by leveraging this deterministic communication behavior and statically setting the congestion window to minimize the chance of deadlock while obtaining the same performance as baseline PFC. Designing such CC and studying the multi-workload cases are the subject of our future studies.

## V. RELATED WORKS

Although congestion control schemes are vastly studied for general datacenter networks, there is little research on CC algorithms for distributed training. Rashidi et al. [18] studied different TCP network configurations for distributed training of recommendation models. Moreover, recent works have shown the applicability of reinforcement learning for internet network congestion control [39][40]. In contrast, this paper mainly focuses on RoCE CC algorithms for distributed training algorithms.

## VI. CONCLUSION

In this paper, we provided a detailed analysis on comparing different proposed RoCE congestion control schemes for distributed training workload/platforms. We highlighted the fact that distributed training workloads/platforms have distinct characteristics that set them apart compared to the general datacenter environments. Therefore, it is essential to tune the networking stack for such platforms separately. We compared the baseline PFC against DCQCN, DCTCP, Timely, and HPCC congestion control schemes. We presented their behavior for both microbenchmark workloads and end-to-end distributed training workloads. In future, we also plan to study the effect of congestion control when multiple training workloads share the same datacenter resources. Based on our analysis, we defined the characteristics of a tuned congestion control scheme for distributed training, which is the subject of our future work.

## REFERENCES

[1] K. He *et al.*, "Deep residual learning for image recognition," 2015.

[2] B. Kazemivash *et al.*, "A novel 5d brain parcellation approach based on spatio-temporal encoding of resting fmri data from deep residual learning," *Journal of Neuroscience Methods*, vol. 369, p. 109478, 2022.

[3] Kazemivash *et al.*, "Bparc: A novel spatio-temporal (4d) data-driven brain parcellation scheme based on deep residual networks," in *IEEE 20th International Conference on Bioinformatics and Bioengineering*, ser. BIBE, 2020, pp. 1071–1076.

[4] "Google Open-Sources Trillion-Parameter AI Language Model Switch Transformer," https://www.infoq.com/news/2021/02/google-trillion-parameter-ai/, 2021.

[5] A. Abdi *et al.*, "Restructuring, pruning, and adjustment of deep models for parallel distributed inference," 2020. [Online]. Available: https://arxiv.org/abs/2008.08289

[6] M. Naumov *et al.*, "Deep learning recommendation model for personalization and recommendation systems," 2019.

[7] B. Klenk *et al.*, "An in-network architecture for accelerating shared-memory multiprocessor collectives," in *47th Annual International Symposium on Computer Architecture*, ser. ISCA, 2020, pp. 996–1009.

[8] A. Sapio *et al.*, "Scaling distributed machine learning with in-network aggregation," *CoRR*, vol. abs/1903.06701, 2019.

[9] M. Cho *et al.*, "Blueconnect: Decomposing all-reduce for deep learning on heterogeneous network hierarchy," *IBM Journal of Research and Development*, vol. 63, no. 6, pp. 1:1–11, Oct. 2019.

[10] S. Rashidi *et al.*, "Themis: A network bandwidth-aware collective scheduling policy for distributed training of dl models," in *Proceedings of the 49th Annual International Symposium on Computer Architecture*, ser. ISCA '22. New York, NY, USA: Association for Computing Machinery, 2022, p. 581–596. [Online]. Available: https://doi.org/10.1145/3470496.3527382

[11] "Cloud TPU," https://cloud.google.com/tpu, 2018.

[12] D. Mudigere *et al.*, "Software-hardware co-design for fast and scalable training of deep learning recommendation models," 2021.

[13] "Microsoft Brainwave," https://www.microsoft.com/en-us/research/project/project-brainwave/, 2022.

[14] "NVIDIA H100 Tensor Core GPU," https://www.nvidia.com/en-us/data-center/h100/, 2022.

[15] "Gaudi Training Platfrom White Paper," https://habana.ai/wp-content/uploads/2019/06/Habana-Gaudi-Training-Platform-whitepaper.pdf, 2019.

[16] "NVIDIA DGX-2," 2019. [Online]. Available: https://www.nvidia.com/en-us/data-center/dgx-2/

[17] A. Shah *et al.*, "Synthesizing collective communication algorithms for heterogeneous networks with TACCL," 2021.

[18] S. Rashidi *et al.*, "Scalable distributed training of recommendation models: An astra-sim + ns3 case-study with tcp/ip transport," in *2020 IEEE Symposium on High-Performance Interconnects*, ser. HOTI, 2020, pp. 33–42. [Online]. Available: https://doi.org/10.1109/HOTI51249.2020.00020

[19] C. Guo *et al.*, "Rdma over commodity ethernet at scale," in *Proceedings of the 2016 ACM SIGCOMM Conference*, ser. SIGCOMM '16. New York, NY, USA: Association for Computing Machinery, 2016, p. 202–215.

[20] Y. Zhu *et al.*, "Congestion control for large-scale rdma deployments," *SIGCOMM Comput. Commun. Rev.*, vol. 45, no. 4, p. 523–536, aug 2015.

[21] Y. Le *et al.*, "Rogue: Rdma over generic unconverged ethernet," in *ACM Symposium on Cloud Computing*, ser. SoCC '18. New York, NY, USA: Association for Computing Machinery, 2018, p. 225–236.

[22] R. Mittal *et al.*, "Timely: Rtt-based congestion control for the datacenter," in *Proceedings of the 2015 ACM Conference on Special Interest Group on Data Communication*, ser. SIGCOMM '15. New York, NY, USA: Association for Computing Machinery, 2015, p. 537–550.

[23] Y. Li *et al.*, "Hpcc: High precision congestion control," in *Proceedings of the ACM Special Interest Group on Data Communication*, ser. SIGCOMM '19. New York, NY, USA: Association for Computing Machinery, 2019, p. 44–58.

[24] S. Rashidi *et al.*, "Astra-sim: Enabling sw/hw co-design exploration for distributed dl training platforms," in *2020 IEEE International Symposium on Performance Analysis of Systems and Software*, ser. ISPASS, 2020, pp. 81–92. [Online]. Available: https://doi.org/10.1109/ISPASS48437.2020.00018

[25] "High-Precision-Congestion-Control," https://github.com/alibaba-edu/High-Precision-Congestion-Control, 2020.

[26] W. Won *et al.*, "Exploring multi-dimensional hierarchical network topologies for efficient distributed training of trillion parameter dl models," 2021. [Online]. Available: https://arxiv.org/abs/2109.11762

[27] R. Thakur *et al.*, "Optimization of collective communication operations in mpich," *International Journal of High Performance Computing Applications*, vol. 19, no. 1, pp. 49–66, Feb. 2005.

[28] G. Wang *et al.*, "Blink: Fast and generic collectives for distributed ml," in *2020 Machine Learning and Systems*, ser. MLSys, 2020, pp. 172–186.

[29] J. Dong *et al.*, "Eflops: Algorithm and system co-design for a high performance distributed training platform," in *2020 IEEE International Symposium on High Performance Computer Architecture*, ser. HPCA, 2020, pp. 610–622.

[30] "NVIDIA Collective Communication Library (NCCL)," 2017. [Online]. Available: https://developer.nvidia.com/nccl

[31] Intel, "Intel oneAPI Collective Communications Library," 2020. [Online]. Available: https://software.intel.com/content/www/us/en/develop/tools/oneapi/components/oneccl.html

[32] P. Devkota *et al.*, "Performance of quantized congestion notification in tcp incast scenarios of data centers," in *2010 IEEE International Symposium on Modeling, Analysis and Simulation of Computer and Telecommunication Systems*, 2010, pp. 235–243.

[33] M. Alizadeh *et al.*, "Data center tcp (dctcp)," in *Proceedings of the ACM SIGCOMM 2010 Conference*, ser. SIGCOMM '10. New York, NY, USA: Association for Computing Machinery, 2010, p. 63–74. [Online]. Available: https://doi.org/10.1145/1851182.1851192

[34] R. B. Basat *et al.*, "Pint: Probabilistic in-band network telemetry," 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*, ser. SIGCOMM '20. New York, NY, USA: Association for Computing Machinery, 2020, p. 662–680.

[35] "ASTRA-SIM: Enabling SW/HW Co-Design Exploration for Distributed DL Training Platforms," https://github.com/astra-sim/astra-sim.git, 2020.

[36] S. Rashidi *et al.*, "Enabling compute-communication overlap in distributed deep learning training platforms," in *2021 ACM/IEEE 48th Annual International Symposium on Computer Architecture (ISCA)*, 2021, pp. 540–553. [Online]. Available: https://doi.org/10.1109/ISCA52012.2021.00049

[37] M. Cowan *et al.*, "Gc3: An optimizing compiler for gpu collective communication," 2022.

[38] "ConnectX SmartNICs," https://www.nvidia.com/en-in/networking/ethernet-adapters/, 2021.

[39] N. Jay *et al.*, "Internet congestion control via deep reinforcement learning," *CoRR*, vol. abs/1810.03259, 2018.

[40] H. Jiang *et al.*, "When machine learning meets congestion control: A survey and comparison," *Computer Networks*, vol. 192, p. 108033, 2021.