

# TACOS: Topology-Aware Collective Algorithm Synthesizer for Distributed Machine Learning

William Won\*, Midhilesh Elavazhagan†, Sudarshan Srinivasan†, Swati Gupta‡, and Tushar Krishna\*

\*Georgia Institute of Technology †Intel ‡Massachusetts Institute of Technology

\*william.won@gatech.edu †midhilesh.elavazhagan@intel.com ‡sudarshan.srinivasan@intel.com

‡swatig@mit.edu \*tushar@ece.gatech.edu

**Abstract**—The surge of artificial intelligence, particularly large language models, has driven the rapid development of large-scale machine learning clusters. Executing distributed models on these clusters is often constrained by communication overhead, making efficient utilization of available network resources crucial. As a result, the routing algorithm employed for collective communications (i.e., collective algorithms) plays a pivotal role in determining overall performance. Unfortunately, existing collective communication libraries for distributed machine learning are limited by a fixed set of basic collective algorithms. This limitation hinders communication optimization, especially in modern clusters with heterogeneous and asymmetric topologies. Furthermore, manually designing collective algorithms for all possible combinations of network topologies and collective patterns requires heavy engineering and validation efforts. To address these challenges, this paper presents TACOS, an autonomous synthesizer capable of automatically generating topology-aware collective algorithms tailored to specific collective patterns and network topologies. TACOS is highly flexible, synthesizing an All-Reduce algorithm for a heterogeneous 128-NPU system in just 1.08 seconds, while achieving up to a  $4.27\times$  performance improvement over state-of-the-art synthesizers. Additionally, TACOS demonstrates better scalability with polynomial synthesis times, in contrast to NP-hard approaches which only scale to systems with tens of NPUs. TACOS can synthesize for 40K NPUs in just 2.52 hours.

**Index Terms**—collective communication, collective algorithm synthesizer, distributed machine learning

## I. INTRODUCTION

The advancement of Artificial Intelligence (AI) has led to a recent surge in the design of specialized distributed High-performance Computing (HPC) platforms tailored for AI, often called AI supercomputers. Examples include Intel Gaudi [1], Google Cloud TPU [2], NVIDIA HGX [3], Cerebras Andromeda [4], Tesla Dojo [5], and many more. From a bird’s eye view, these platforms incorporate multiple Neural Processing Units (NPUs, such as GPUs, TPUs, or ASICs) at the endpoints, interconnected by custom high-speed fabrics. These AI supercomputers are of particular importance for enormous Deep Neural Network (DNN) models, such as Large Language Models (LLMs), due to their substantial compute and memory demands [6].

It is important to note that different AI clusters employ diverse network topologies. For instance, point-to-point connections are used in Intel Gaudi [1] and Google TPU [2], hierarchical switches are utilized in NVIDIA HGX [7], mesh architectures are employed in Tesla Dojo [5] and Cerebras Andromeda [4], and DragonFly [8] connections are used in

TABLE I  
ALL-REDUCE ALGORITHMS AND THEIR TARGET TOPOLOGIES. TOPOLOGY ABBREVIATIONS DENOTE RI (RING), FC (FULLYCONNECTED), SW (SWITCH), HO (HOMOGENEOUS), AND HT (HETEROGENEOUS).  $\Delta$  DENOTES THE COLLECTIVE PERFORMANCE IS CLOSELY TIED TO SPECIFIC NETWORK CONFIGURATIONS.

Topology-aware All-Reduce Algorithm	Preferred Physical Topology							
	Uni-dimensional			Multi-dim.		Asym-metric		Any
	RI	FC	SW	Ho	Ht	Ho	Ht	Any
Ring [21]	✓							
Direct [22]		✓						
RHD [23]			✓					
DBT [24]			✓					
BlueConnect [25]	✓	✓	✓	$\Delta$	$\Delta$			
Themis [18]	✓	✓	✓	✓	✓			
TTO [26]				$\Delta$		$\Delta$		
C-Cube [27]		$\Delta$		$\Delta$		$\Delta$		
TACOS (this work)	✓	✓	✓	✓	✓	✓	✓	✓

TABLE II  
QUALITATIVE COMPARISON OF COLLECTIVE ALGORITHM SYNTHESIZERS. TACOL IS MARKED  $\Delta$  AS IT REQUIRES ASSUMPTIONS (COMMUNICATION SKETCH) TO MODEL NETWORK HETEROGENEITY AND CONGESTION.

Framework	Network		Synthesis Mechanism		
	Asym-metric	Hetero-geneous	Auton-omous	Removes Conges-tion	Scal-able
SCCL [20]			✓		
Blink [28]	✓		✓		
MultiTree [29]	✓		✓	✓	
TACOL [19]		$\Delta$	$\Delta$	$\Delta$	
TACOS (this work)	✓	✓	✓	✓	✓

Groq [9]. Moreover, these systems often incorporate multiple link technologies such as XeLink [10], NVLink [11], InfiniBand [12], or even optical networks [2]. In summary, topologies in Machine Learning (ML) clusters commonly exhibit *asymmetric*<sup>1</sup> shapes and *heterogeneous*<sup>2</sup> link bandwidths. Considering that communication among NPUs is the main bottleneck of distributed ML [13]–[15], orchestrating communication within these AI platforms remains an open research problem [16]–[20].

Communication within AI workloads presents unique characteristics compared to traditional cloud workloads. In the

<sup>1</sup>For e.g., NPUs at the center vs. edge of a 2D Mesh have different degrees.

<sup>2</sup>All links within a *homogeneous* topology have the same bandwidth.

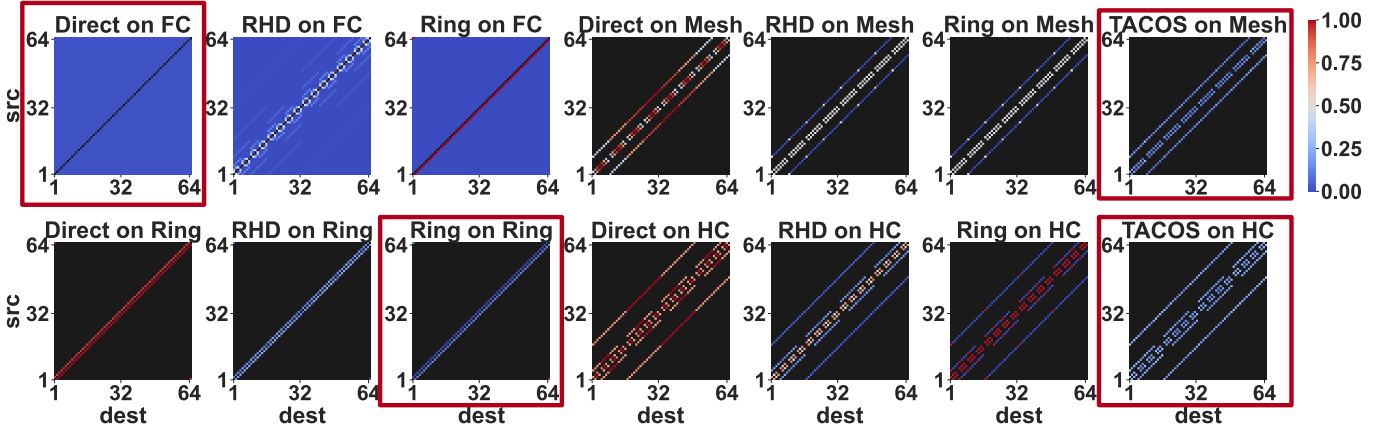


Fig. 1. Heat map of total message size transferred over each link, when running 1 GB All-Reduce using basic algorithms (Direct, RHD, Ring, and TACOS for comparison) over different network topologies (FullyConnected (FC), Ring, 2D Mesh, and 3D Hypercube (HC)). Each cell at  $(src, dest)$  denotes a link connecting NPU  $src$  to NPU  $dest$ . If there is no such link in the topology, the cell is marked black. All values are normalized to the largest value per topology. For every topology, scenarios running topology-aware collective algorithms are marked with a red box, which results in balanced, lowest overall loads (i.e., a cooler heat map).

latter, there can be heavy variability in the traffic pattern and volumes. This has naturally led to a set of past works on dynamically managing network congestion in data center fabrics [30]. In contrast, for ML models, both the traffic pattern and traffic volume are *deterministic* once the workload has been partitioned and mapped [31]. The traffic pattern is *collective* in nature (i.e., one-to-many and many-to-one). Given these characteristics, AI workload communication today is often routed *statically* via Collective Communication Libraries (CCLs) [32]–[35]. These CCLs encompass a collection of predefined collective algorithms known as *topology-aware collectives*. As the name suggests, these algorithms are tailored for operation on specific network topologies [36]. For example, Ring (RI) [21], Direct (DI) [22], Recursive Halving-Doubling (RHD) [23], and Double Binary Tree (DBT) [24] are distinct All-Reduce collective algorithms tailored for fundamental network topologies (a.k.a. basic collective algorithms). Table I exemplifies several predefined collective algorithms and their target topologies.

Naively deploying topology-aware collectives over non-preferred physical topologies could induce link congestions and underutilization of the network resources. This is demonstrated in Fig. 1. When basic collective algorithms are executed over their non-preferred topologies, significant link underutilization and oversubscriptions occur.<sup>3</sup> This is the reason CCLs today have heuristics built-in to switch among basic collective algorithms. Unfortunately, beyond a limited set of basic collective algorithms, there are no known solutions for arbitrary topologies. Given this complexity, when given an unknown topology, CCLs today default to employing the Ring algorithm.<sup>4</sup> However, this approach may result in oversubscription and underutilization of certain links.

To address the aforementioned challenge, recent studies

<sup>3</sup>We use *bidirectional* Ring algorithm and *bidirectional* Ring topology for all the motivation and evaluation data in this paper.

<sup>4</sup>Either one logical ring is mapped over the physical topology, or multiple parallel rings [18], [25]. But the problem of over/undersubscription remains.

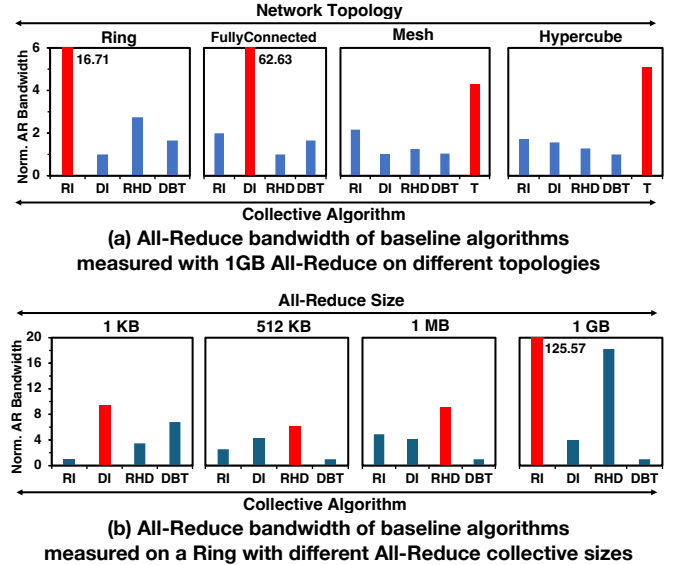
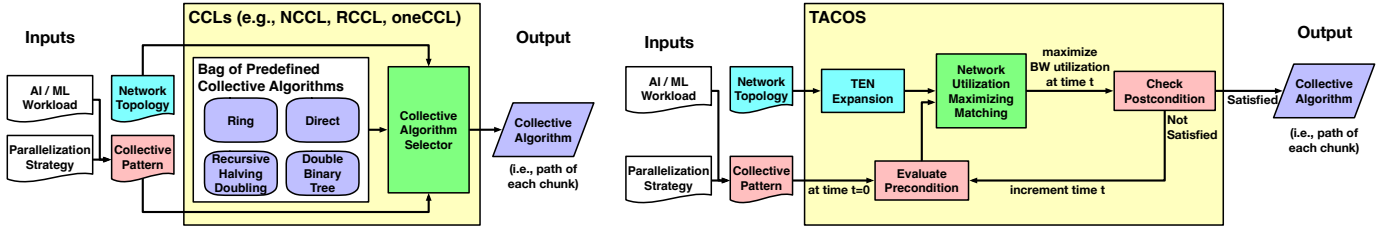


Fig. 2. (a) All-Reduce bandwidth (i.e., collective size  $\div$  collective time) over distinct network topologies with 64 NPUs (using link  $\alpha=0.5\mu s$ ,  $1/\beta=50GB/s$ , as explained in Sec. IV-G), measured with 1 GB collective. For 2D Mesh and 3D Hypercube, we also run the topology-aware algorithm synthesized by TACOS (T). (b) All-Reduce bandwidth for different collective algorithms on a 128-NPU physical Ring (link  $\alpha=30ns$ ,  $1/\beta=150GB/s$ ), measured with varying collective sizes. Results are normalized per each graph by the smallest number.

have demonstrated the feasibility and effectiveness of *autonomously synthesizing* topology-aware collective algorithms when provided with a target physical network and collective patterns [19], [20], [28], [29]. These efforts are outlined in Table II. However, they are limited to supporting a specific subset of topologies, typically homogeneous or symmetric networks. Furthermore, they often overlook network congestion effects during the search, resulting in suboptimal synthesis results. A number of works treat collective synthesis as a global optimization problem, which poses NP-hard complexity to the search space, fundamentally limiting their scalability [37].



(a) Current CCLs: Among predefined collective algorithms, select the best-performing algorithm by estimating collective performance by using the underlying network topology and collective pattern information.

(b) TACOS: Autonomously synthesize a topology-aware collective algorithm by repeatedly executing Network Utilization Maximizing Matching algorithm over a TEN. TACOS thereby synthesizes quality collective algorithm with maximized network utilization.

Fig. 3. (a) Current CCLs execute collectives by selecting an algorithm from a set of predefined implementations. (b) High-level overview of the TACOS framework. The target network topology and collective pattern are provided as inputs. TACOS expands the network into a TEN and evaluates the collective precondition (i.e., which chunks are currently held by each NPU) for time  $t=0$ . Based on this information, TACOS employs a chunk-link matching algorithm that maximizes network resource utilization for the target time  $t=0$ . TACOS iteratively runs this matching process for successive time spans until the postcondition is satisfied (i.e., all NPUs have received every desired chunk). The details are explained in Sec. IV. This procedure yields a topology-aware collective algorithm (i.e., static path of each chunk), which can then be utilized by CCLs in lieu of the predefined topology-unaware basic algorithms.

TABLE III  
COMMON MODEL AND DATA PARALLELIZATION STRATEGIES AND THEIR REQUIRED COLLECTIVE COMMUNICATION PATTERNS.

Parallelization	Reduce-Scatter	All-Gather	All-Reduce
Data Parallelism [22]			✓
Tensor Parallelism [22]			✓
FSDP [41]	✓	✓	
ZeRO [6]	✓	✓	
Hybrid [22]	✓	✓	✓

Consequently, their evaluations have been restricted to only tens of NPUs.

In this work, we propose an autonomous collective synthesizer named **TACOS: Topology-Aware Collective Algorithm Synthesizer**.<sup>5</sup> TACOS architecture is summarized in Fig. 3. Given a network topology and target collective pattern, TACOS autonomously synthesizes a *static* collective algorithm, which can be leveraged by CCLs. This is done by *iteratively maximizing network resource utilization* throughout the course of collective execution. TACOS showcases better quality synthesis results, orders of magnitude better scalability, and encompasses a larger class of topologies compared to previous endeavors. This is enabled by two features. Firstly, we frame collective algorithm synthesis as a *link-chunk matching problem* using a Time-expanded Network (TEN) representation, an acyclic directed graph integrating spatial and temporal dimensions [38]–[40]. Secondly, we employ a *link-chunk matching algorithm* to maximize network utilization and minimize congestion. Here is a summary of our contributions:

- TACOS brings the TEN representation into the domain of distributed ML, which is a widely adopted data structure for traffic and flow optimizations.
- TACOS accommodates a wide range of arbitrary (i.e., heterogeneous and asymmetric) topologies.
- TACOS considers network congestion effects throughout the synthesis process, resulting in high-quality search outcomes.
- TACOS features better scalability, surpassing previous works usually restricted to tens of NPUs. TACOS synthesized collective algorithms for a 40K NPU topology in 2.52 hours.

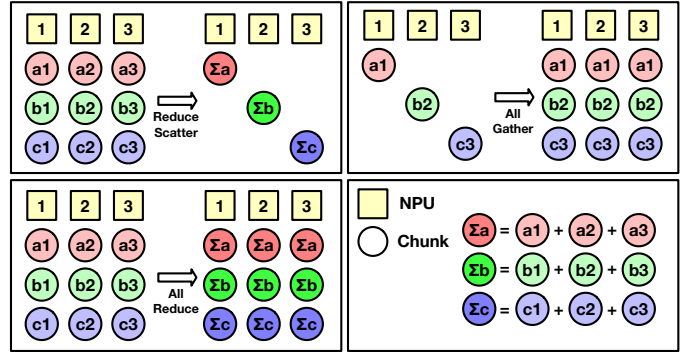


Fig. 4. Common collective patterns used in distributed ML. Each collective communication defines a specific data transfer pattern. For example, in All-Gather, each NPU starts with one chunk (denoted as a circle) and broadcasts it to all other NPUs. In Reduce-Scatter, all NPUs start with  $N$  chunks and end up with one chunk, constructed by summing the corresponding chunks from all other NPUs.

## II. BACKGROUND

### A. Collective Communication Patterns

Communication poses a significant challenge for distributed ML systems due to the dispersal of models (and training data for training tasks) across devices, requiring frequent synchronization among NPUs [18]. These communications can be represented as collective communications [22]. As depicted in Table III, workload parallelization strategies require specific collective patterns to be executed. Fig. 4 illustrates common collective patterns. Each circle represents a chunk, the atomic unit for scheduling collectives across the network. The All-Reduce can be conceptualized as two sequential stages: Reduce-Scatter followed by All-Gather, and it is the most prevalent pattern in synchronous NPU-to-NPU ML executions [42]. To enhance network utilization, a collective can be decomposed into multiple smaller chunks which can be run in parallel [18].

### B. Collective Communication Algorithms

For a given collective communication pattern, there can be several distinct routing algorithms to implement it. These

<sup>5</sup>TACOS is open-sourced at: <https://github.com/astra-sim/tacos>.

are called *collective algorithms* and define the *static path* for each chunk during the execution of collective communication. Example collective algorithms' traffic patterns are depicted in Fig. 5. The intuition for having different algorithms for the same collective pattern is to minimize contention on distinct physical network topologies.

### III. MOTIVATION AND PROBLEM STATEMENT

#### A. Importance of Topology-aware Collectives

The physical characteristics of the network topology (including connectivity, link latencies, and link bandwidths) heavily influence network transmission behavior. We define a *topology-aware collective algorithm* as one intricately optimized alongside the underlying topology to achieve optimal collective performance [23], [43]. Table I summarizes common topology-aware algorithms and their target topologies.

Example traffic patterns of basic collective algorithms are depicted in Fig. 5. When these patterns are mapped over their *preferred* topologies, they can efficiently utilize links in the topology every cycle without any contention. However, executing a collective algorithm on its unpreferred physical topology leads to contention (i.e., oversubscription) on certain links and low or even no utilization (i.e., undersubscription) on others. This effect can be observed in the link subscription heat map in Fig. 1. We also quantify these effects in Fig. 2(a). On the Ring network, the topology-aware Ring algorithm exhibited a  $16.71\times$  higher All-Reduce bandwidth than the Direct algorithm. Conversely, on the FullyConnected topology, the Direct algorithm demonstrated a performance improvement of  $62.63\times$  over the Ring algorithm. Fig. 2(b) also demonstrates that the optimal topology-aware collective could change depending on the target collective pattern. Unlike the 1 GB All-Reduce case, for 1 KB, the Direct algorithm outperformed Ring since each network transmission is latency-bound and prefers short-distance collective algorithms.

These results emphasize the importance of designing optimal topology-aware collective algorithms that encapsulate both communication patterns and network topology. Modern CCLs [32], [33], [35] implement a range of predefined basic collective algorithms, selecting them based on network features as shown in Fig. 3(a). Unfortunately, this fundamentally limits the physical topologies they can efficiently support.

#### B. Heterogeneity, Asymmetry, and Scale of ML Systems

To provide maximum network resources, state-of-the-art distributed ML platforms are leveraging various network topologies and technologies, introducing bandwidth heterogeneity and asymmetry. Notable examples include asymmetric 2D Mesh topology with wafer-scale technologies as implemented by Cerebras CS-2 [4] and Tesla Dojo [5] systems. Heterogeneous multi-dimensional scale-up and scale-out switches in NVIDIA DGX clusters [11], [12], and the integration of scale-up and photonic technologies by Google Cloud TPU [2], [44] are other examples. Moreover, the scale of recent ML clusters is also increasing to accommodate larger workloads, incorporating thousands to even tens of thousands of NPUs.

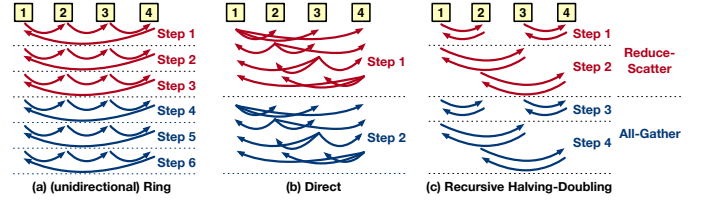


Fig. 5. Example All-Reduce algorithms and their traffic patterns. In Reduce-Scatter, each red arrow represents sending and adding a chunk to the local data. For All-Gather, each blue arrow signifies forwarding a chunk. Depending on the network connectivity, each step may encounter link congestion and underutilization.

This trend is evident in Google Cloud TPUv5 pods with 8,960 TPUs [45], the OpenAI cluster with 7,500 GPUs [46], and Meta's Research Supercluster equipped with 24K GPUs [47].

#### C. Challenge: Designing Topology-aware Collectives

Given the importance of topology-aware collective algorithms in distributed ML, designing them for commonly used networks is an active research field. For example, recent proposals devise All-Reduce algorithms for 2D Mesh [48] and DragonFly [49] topologies. However, this process requires engineering and validation efforts for each new topology variant. Flattened Butterfly [50], MegaFly [51], SlimFly [52], and Tofu [53] are just a few examples that do not yet have specialized collective algorithms and default to baseline collective algorithms, which could lead to network resource underutilization. Even existing collective algorithms are often tuned under multiple assumptions, such as homogeneous link bandwidths or specific message sizes, making them susceptible to potential performance degradation in practice.

#### D. Solution: Topology-aware Collective Synthesizer

To address the challenges inherent in designing topology-aware collective algorithms manually, the development of an *autonomous synthesis framework capable of automatically generating topology-aware collective algorithms* is imperative. Such a toolchain can alleviate the burden on human experts to manually configure the routing of each chunk. Additionally, by automatically exploring the design space, these frameworks have the potential to generate higher-quality collective algorithms for specific targets. This is particularly crucial as network topologies become larger and more intricate, further complicating the task of design and optimization.

## IV. TACOS

#### A. Background on Time-expanded Network

We start introducing TACOS by bringing the notion of Time-expanded Network (TEN) to the distributed ML domain. TEN is widely used in flow optimization problems [39], [54], [55] because it enables judicious consideration of both network topology and timing information in a single data structure.

A visual example of a homogeneous, asymmetric 3-NPU topology is shown in Fig. 6(a), as well as its TEN representation in Fig. 6(b). Each NPU in the base topology forms a column of vertices in the TEN representation, with each



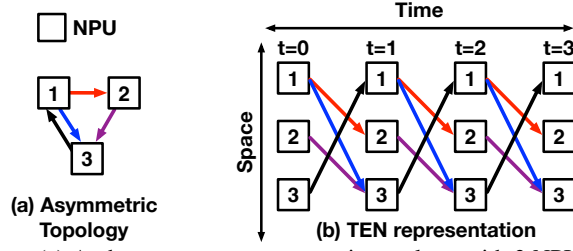
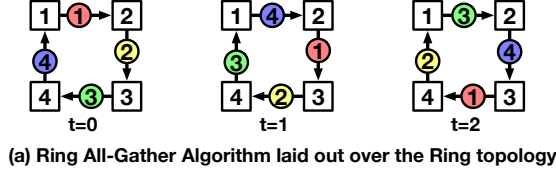
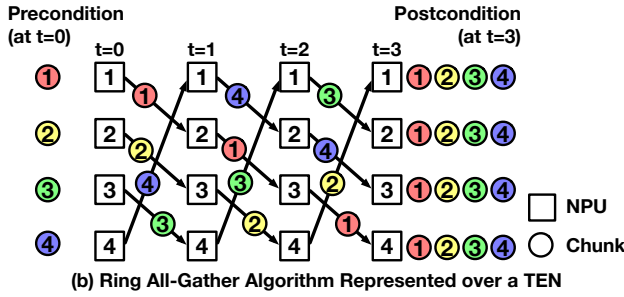


Fig. 6. (a) An homogeneous, asymmetric topology with 3 NPUs and 4 links. (b) TEN expansion of the network up to  $t=3$ .



(a) Ring All-Gather Algorithm laid out over the Ring topology



(b) Ring All-Gather Algorithm Represented over a TEN

Fig. 7. (a) Unidirectional Ring All-Gather algorithm laid out over a unidirectional Ring topology across different time spans. (b) Algorithm shown in (a) captured using the TEN-based representation.

column corresponding to a unique time span. This column is then replicated across a time range (up to  $t=3$  in this example). The topology exhibits four unidirectional connections:  $1 \rightarrow 2$ ,  $1 \rightarrow 3$ ,  $2 \rightarrow 3$ , and  $3 \rightarrow 1$ . These connections translate into distinct edges in the TEN representation, linking different time spans flowing from left to right.

This TEN representation effectively integrates the temporal and spatial information of the network communication behavior into a single data structure. Consider the scenario where a chunk starts being transmitted from NPU 1 to NPU 2 at  $t=1$ . It can be seen from the TEN that this chunk will reach its destination at  $t=2$ . Similarly, the absence of a direct link from NPU 3 to NPU 2 in the TEN illustrates that such a direct transmission is impossible due to the lack of the physical link in the network.

### B. Representing Collective Algorithms using TEN

TEN offers a concise mechanism to represent collective algorithms in a single data structure. One example is shown in Fig. 7. Depicted in Fig. 7(a) is an unidirectional Ring All-Gather algorithm laid over the Ring topology across distinct time spans. The same collective algorithm is represented over the TEN representation in Fig. 7(b). The leftmost segment of the TEN depicts the *precondition* (i.e., chunks initially held by each NPU) at  $t=0$ , while the rightmost segment represents the *postcondition* (i.e., chunks held at the end) at  $t=3$ .

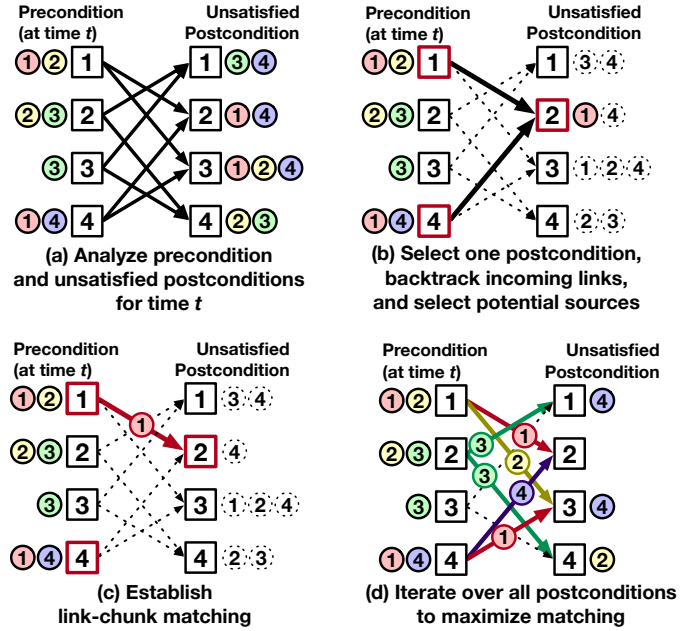


Fig. 8. Network Utilization Maximizing Matching algorithm. For a given time span  $t$  of a TEN: (a) The pre/postconditions of a collective communication are evaluated. (b) An unsatisfied postcondition is randomly selected, and the destination NPU is backtracked over the TEN to identify candidate source NPUs capable of providing the requested chunk. (c) One candidate is randomly chosen and a link-chunk match is made. (d) Steps (b)-(c) are repeated exhaustively over all unsatisfied postconditions to maximize the number of link-chunk matches, optimizing link utilization for the given time  $t$ .

In this visualization, a chunk transmission between two NPUs is symbolized by the chunk occupying a TEN link, which we term a *link-chunk match*. For this example, all links in the TEN structure correspond to matched chunks, indicating maximized resource utilization throughout the course of the collective execution. Furthermore, *each link matches up to one chunk*, eliminating network contention from multiple payloads assigned over a single link at a time.

To the best of our knowledge, this work is the first to introduce TEN to the domain of distributed ML and to represent collective algorithms atop it. We believe this formulation has value beyond the specific collective synthesis mechanism presented in this paper.

### C. TACOS' Approach to Collective Algorithm Synthesis

The TEN-based representation employed by TACOS enables it to frame the collective synthesis problem as a link-chunk match-making challenge, rather than as a global optimization problem. This strategy offers a significant advantage, as solving global optimizations (e.g., Integer Linear Programming (ILP) formulations) is an NP-hard problem [37].

Instead, in order to synthesize a collective algorithm, TACOS aims to *establish and maximize the number of link-chunk matches* for individual chunks onto the TEN links. To synthesize the collective algorithm, TACOS first constructs a TEN for  $t=0$  for the target network topology. It then proceeds to *execute a link-chunk matching algorithm* (described in Sec. IV-D)

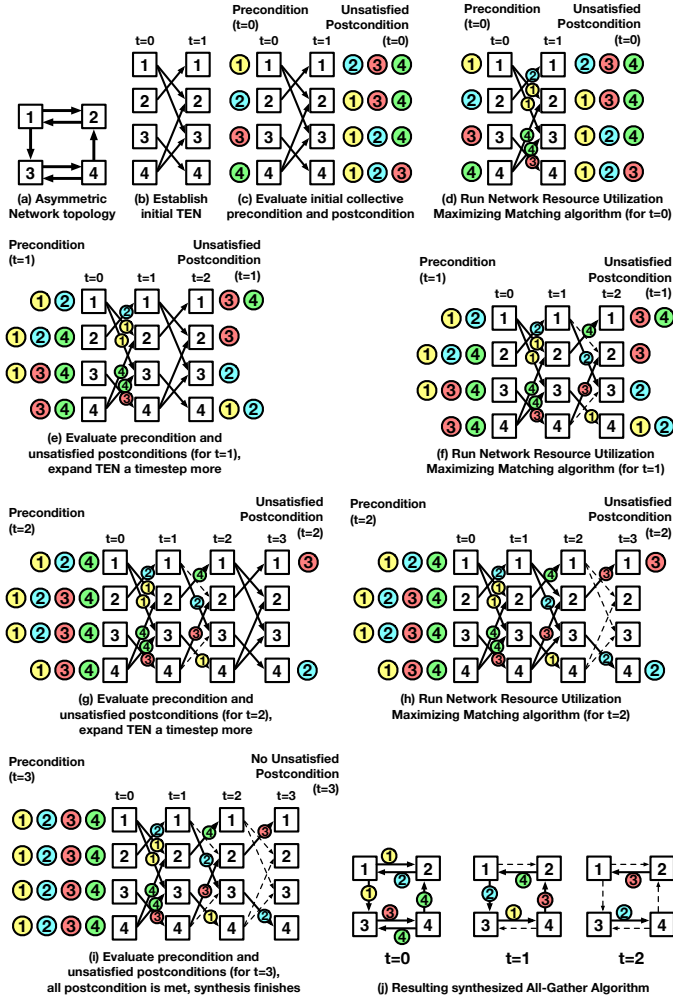


Fig. 9. End-to-end overview of how TACOS constructs an All-Gather algorithm for a homogeneous, asymmetric 4-NPU topology.

to maximize network utilization at  $t=0$ . The *TEN* is then expanded by one more time span, and the process repeats until all postconditions are met.

This iterative synthesis approach maximizes network utilization over the entire course of the collective execution, thereby producing a high-performance collective algorithm. Moreover, the framework *inherently mitigates link congestion* by allowing each *TEN* link to accommodate only one chunk. Importantly, TACOS supports arbitrary networks since any heterogeneous or asymmetric topologies can still be represented in *TEN*, eliminating the need for specific homogeneity or symmetry assumptions.

#### D. Network Utilization Maximizing Matching Algorithm

TACOS employs a swift *network utilization maximizing matching algorithm* to find link-chunk matches for a given time span  $t$ . This algorithm is summarized in Fig. 8. The matching process begins in Fig. 8(a), where TACOS evaluates the preconditions and unsatisfied postconditions of each NPU at the given time span  $t$  to determine the chunks requested to be received by each NPU. TACOS then randomly selects one unsatisfied postcondition to evaluate potential link-chunk

#### Algorithm 1 Utilization Maximizing Matching Algorithm

**Condition:** Provided *TEN* at time  $t$

```

// Get unsatisfied postconds
for  $d \leftarrow \forall \text{NPUs}$  do
  for  $c \leftarrow \forall \text{Chunks}$  do
    if  $\text{postcond}[d][c] \ \&\& \ !\text{precond}[d][c]$  then
      Add  $(d, c)$  to unsatisfied_postcond

// Make link-chunk matches
random_shuffle(unsatisfied_postcond)
for  $(d, c) \leftarrow \forall \text{unsatisfied\_postcond}$  do
  // Select  $s$  NPU that can provide  $c$  to  $d$ 
   $s \leftarrow \text{backtrack}(\text{TEN}[t][*][d])$ 
   $\text{candidate\_srcs} \leftarrow \text{List}(\forall s: \text{precond}[s][c] \text{ is True})$ 
   $\text{chosen\_src} \leftarrow \text{random\_select}(\text{candidate\_srcs})$ 

// Mark link-chunk match
 $\text{precond}[d][c] \leftarrow \text{mark True}$ 
 $\text{TEN}[t][\text{chosen\_src}][d] \leftarrow \text{mark } c$ 

```

#### Algorithm 2 TACOS End-to-End Synthesis

**Condition:** Provided target topology and collective

**Initialize:**  $\text{TEN}[t = 0]$ , pre/postconds

```

while  $\exists \text{unsatisfied\_postconds}$  do
  // Expand TEN by 1 time span and run matching
   $t \leftarrow t + 1$ 
  Expand  $\text{TEN}[t]$ 
  Utilization_Maximize_Matching( $\text{TEN}[t]$ ) (Alg. 1)

```

matches. For instance, in Fig. 8(b), NPU 2 and chunk 1 are selected. The algorithm then backtracks the *TEN* graph to identify a set of potential source NPUs capable of supplying the requested chunk, as represented by the bold arrows in the figure. These source NPUs are assessed to determine which NPUs currently possess and can provide the desired chunk. In the illustrated scenario, both NPU 1 and NPU 4 fulfill these criteria. From this pool of candidates, one source NPU is randomly selected, resulting in a successful link-chunk match as depicted in Fig. 8(c). In this example, NPU 1 is chosen to cater to chunk 1, and the corresponding *TEN* link (NPU 1 to 2) becomes occupied by the established match. This sequence repeats for all unsatisfied postconditions, yielding the outcome in Fig. 8(d). The algorithm is summarized in Alg. 1.

For a given time span  $t$ , the matching algorithm exhaustively iterates over all unsatisfied postconditions and makes random selections whenever possible. Therefore, the algorithm tries to *maximize the number of successful link-chunk matches*, thereby ensuring nearly maximal utilization of network link resources at the given time  $t$ . Additionally, as only one chunk can be matched over a link, this algorithm effectively eliminates link congestion.

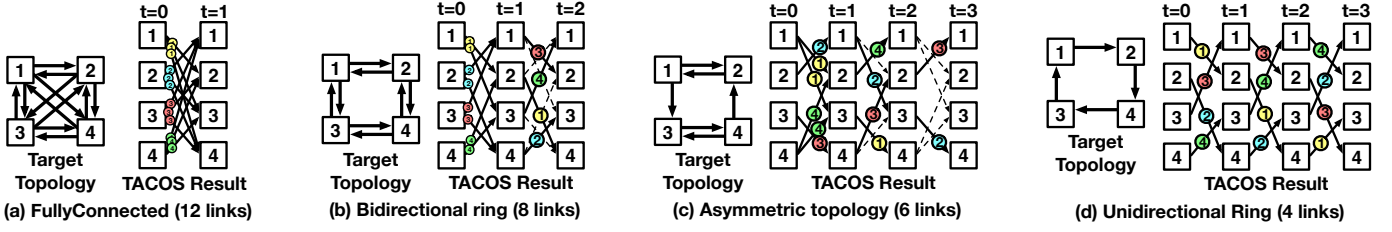


Fig. 10. Distinct target topologies with 4 NPUs but different numbers of links, and their All-Gather synthesis results over TEN using TACOS. As the network connection becomes more sparse, TACOS is required to expand the TEN for more time spans to satisfy all collective postconditions. Still, note that the matching algorithm was able to maximize the network utilization for every time span.

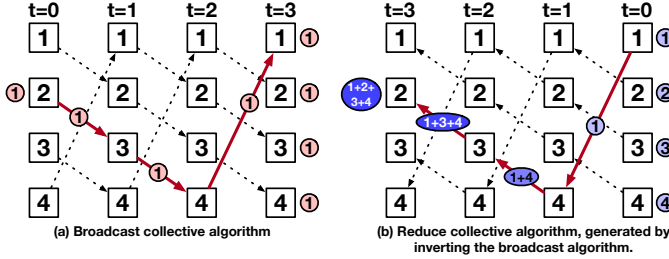


Fig. 11. Synthesis of collectives with reduction operations (e.g., Reduce-Scatter). (a) TACOS first constructs a TEN whose  $(src, dest)$  are reversed to  $(dest, src)$  (reversed link directions) and synthesizes corresponding non-reduction collectives. (b) Then, TACOS reverses the synthesized collectives back to get the final combining collectives.

### E. End-to-end TACOS Collective Synthesis

To help understand the end-to-end synthesis process, Fig. 9 provides an illustrative example synthesizing an All-Gather algorithm for a 4-NPU asymmetric network. Fig. 9(a) depicts the target topology, while Fig. 9(b) shows the initial TEN representation of the network for  $t=0$ . First, preconditions and postconditions are evaluated for  $t=0$ , as shown in Fig. 9(c). Subsequently, the link-chunk matching algorithm is executed, resulting in maximized network utilization for  $t=0$  as demonstrated in Fig. 9(d). The preconditions and postconditions are updated for  $t=1$ , and since there are remaining unsatisfied postconditions, the TEN is expanded for one more time span, resulting in the structure depicted in Fig. 9(e). This process repeats iteratively: the link-chunk matching algorithm is executed for  $t=1$ , once again maximizing network utilization, as shown in Fig. 9(f). The pre/postcondition re-evaluation, TEN expansion, and link-chunk matching are repeated for  $t=2$ , as depicted in Fig. 9(g-i). Once all postconditions are satisfied, the synthesis terminates at Fig. 9(i). The resultant collective algorithm, laid out on the original network topology, is summarized in Fig. 9(j), ready to be deployed by CCLs.

This end-to-end synthesis algorithm is summarized in Alg. 2. We first initialize TEN and postconditions for  $t = 0$ . Until all postconditions are met, TACOS simply expands the TEN by one more time span and runs the utilization maximization matching process. Since TACOS only schedules communications between neighboring NPUs, the TACOS-synthesized algorithm is deadlock-free [29].

Fig. 10 illustrates how TACOS synthesis results react to changes in the underlying topology. All target topologies in Fig. 10(a-d) have 4 NPUs but different numbers of physical

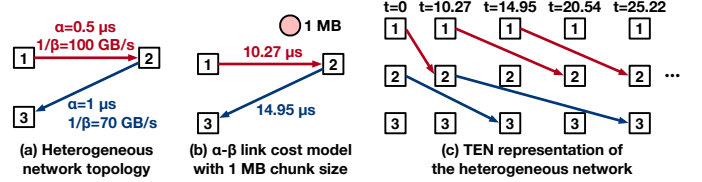


Fig. 12. (a) A heterogeneous 3-NPU topology. (b)  $\alpha$ - $\beta$  model applied to calculate link costs. (c) TEN representation of the heterogeneous topology, expanded up to  $t=25.22 \mu s$ .

links connecting them. For the FullyConnected topology, all postconditions can be satisfied in a single shot, effectively yielding a Direct algorithm. As connectivity becomes scarcer, the execution of the collective takes longer. It is important to note, however, that TACOS remains effective in maximizing link utilization for each discrete time span, resulting in optimal collective algorithms with minimized collective time.

**Collectives with Reduction Operations.** Certain collectives, such as Reduce-Scatter or All-Reduce, involve reduction operations among chunks. TACOS handles such collectives by synthesizing their non-combining counterpart collective patterns. An illustrative example with Broadcast and Reduce collectives is shown in Fig. 11. To synthesize Reduce collective, all  $(src, dest)$  links are inverted into  $(dest, src)$  and the corresponding Broadcast algorithm gets synthesized. By reversing the synthesized algorithm back, TACOS can obtain the Reduce algorithm. Similarly, search for a Reduce-Scatter algorithm can be accomplished by synthesizing an All-Gather algorithm and reversing the sender-receiver of all communication pairs. An All-Reduce operation can be synthesized by combining a Reduce-Scatter followed by an All-Gather.

### F. TACOS for Heterogeneous Networks

The basics of conducting link-chunk matching over TEN remain consistent even for heterogeneous networks.

**Heterogeneous Topology in TEN.** To delineate heterogeneous link costs, we adopt the  $\alpha$ - $\beta$  model [18]–[20], [56], [57].  $\alpha$  characterizes the latency of the link, whereas  $\beta$  represents the reciprocal of link bandwidth (i.e., serialization delay per unit message size). Given a chunk size  $n$ , we can translate each link transmission cost into a single number,  $\alpha + \beta \times n$ . Once the link cost is set, a TEN representation can be drawn. Fig. 12 shows this process. The network configuration in Fig. 12(a) features a heterogeneous topology. Each link delay is derived using the  $\alpha$ - $\beta$  model in Fig. 12(b) (assuming a 1 MB chunk in this scenario). Fig. 12(c) shows the resultant

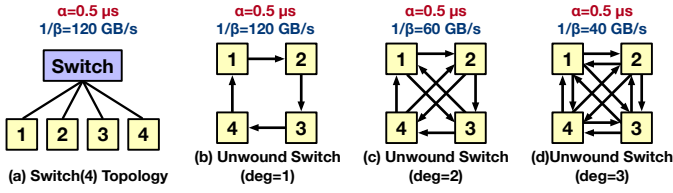


Fig. 13. (a) A Switch fabric with 4 NPUs. (b)-(d) Unwinding the switch network with degrees 1 through 3.

TEN representation, which can then be absorbed by the TACOS framework. Although we have not considered NPU or memory performance during synthesis in this work, these factors can be also incorporated into the  $\alpha$  and  $\beta$  costs of each link.

**Prioritizing Lower-cost Links.** When conducting link-chunk matching, multiple links with different transmission costs may emerge as potential matches due to network heterogeneity. To minimize collective time, TACOS prioritizes the link with the lowest cost when making link-chunk matches.

#### G. TACOS for Switch-based Networks

A Switch topology offers versatile NPU connectivity, but unregulated use leads to sub-optimal performance due to network contentions or even deadlocks. To address this, we propose a mechanism to unfold Switch networks into fixed, point-to-point connections. With a degree  $d$  unwinding approach for an  $N$ -NPU Switch, we establish  $d$  outgoing links from each NPU  $i$ , connecting to NPUs  $(i+1)$ ,  $(i+2)$ ,  $\dots$ ,  $(i+d) \pmod{N}$ .<sup>6</sup> During this transition, the cost  $\alpha$  associated with each point-to-point link remains consistent. However, due to shared link bandwidth, the  $\beta$  cost multiplies by  $d \times$ . Fig. 13(a) illustrates a 4-NPU switch while Fig. 13(b-d) depict its corresponding  $d=1-3$  unwinding. Once the switch is unwound, the topology can be represented in TEN and used by TACOS.

### V. METHODOLOGY

#### A. Baseline Collective Algorithms and Synthesizers

For performance evaluations, we compare TACOS against the following baseline All-Reduce algorithms.

- **Ring, Direct.** These are basic collective algorithms showing traffic patterns depicted in Fig. 5. In particular, Ring currently serves as the default algorithm across most CCLs [32], [33], [35].
- **RHD, DBT.** We also assessed RHD [23] and DBT [24] in Table V since they are only suited for networks with a power-of-two number of NPUs.
- **BlueConnect, Themis, C-Cube.** These are topology-aware collective algorithms manually crafted to accommodate specific network topologies. BlueConnect [25] develops a multi-rail All-Reduce algorithm for a symmetric hierarchical network. Themis [18] further optimizes BlueConnect

<sup>6</sup>This unwinding scheme neither covers all potential point-to-point configurations nor dynamic Switch connectivity across TEN time spans. Notably,  $d=1$  and  $d=N-1$  unwinding are suitable for bandwidth and latency-critical collective synthesis, respectively [19]. However, developing a more flexible Switch unwinding methodology remains a subject for future investigation.

TABLE IV  
TOPOLOGIES EVALUATED IN THIS PAPER.

Topology	Heterogenous	Asymmetric
Ring (RI)		
FullyConnected (FC)		
2D Torus		
3D Torus		
2D Mesh (Mesh)		✓
3D Hypercube (HC)		✓
2D Switch	✓	
3D Ring-FC-Switch (3D-RFS)	✓	
DragonFly (DF)	✓	✓

through improved chunk-level scheduling. Meanwhile, C-Cube [27] manually maps two binary trees on DGX-1 and executes two tree-based All-Reduce algorithms in parallel.

- **MultiTree, TACCL.** MultiTree [29] synthesizes collective algorithms for homogeneous networks via spanning tree construction, while TACCL [19] employs an ILP-based approach for symmetric heterogeneous topologies.<sup>7</sup>
- **Ideal.** Finally, to show TACOS' absolute synthesis quality, we also augmented all results with theoretically ideal collective performance. This can be derived from the topology diameter (the minimum latency for the farthest two NPUs to communicate) for  $\alpha$  costs, and the bottleneck serialization delay (injection and ejection time of all chunks) for  $\beta$ .

$$\text{Ideal} = \frac{\text{CollectiveSize} \times 2(n-1)/n}{\min_{N \in \mathcal{V}_{\text{NPU}}} (\text{BW}_N)} + \text{Diameter}$$

#### B. Target Topologies and Synthesis Environment

Table IV lists the topologies studied in this work, covering both homogeneous and heterogeneous networks, as well as both symmetric and asymmetric topologies.<sup>8</sup> TACOS synthesis time is measured using the Intel Xeon E5-2699v3 CPU.

#### C. Simulation Infrastructure

For this work, a simulator was necessary to study performance across the diverse spectrum of topologies that we examined (Table IV). We use the ASTRA-sim distributed ML simulator [22], [58] to evaluate the performance of the baseline collective algorithms and their corresponding impact on full workloads. ASTRA-sim includes NCCL-validated implementations of the basic collective algorithms. It also models chunking (i.e., breaking a large collective into smaller chunks) and its related effects in BlueConnect and Themis.

**Network Simulation Backend.** To model the network over-subscription and congestion behaviors incurred by topology-unaware collectives, while enabling the simulation of large-scale network topologies, we enhanced the network modeling of ASTRA-sim by creating and attaching a congestion-aware analytical network simulation backend. The congestion-aware analytical backend simulates a message transfer by simulating

<sup>7</sup>C-Cube and MultiTree implementations are not public, and we used their reported numbers in the paper for an apples-to-apples comparison. Since TACCL offers limited topology options, we implemented a TACCL-like baseline by integrating its ILP formulation over our TEN representation.

<sup>8</sup>Unless otherwise specified, we set link  $\alpha=0.5 \mu s$  and  $1/\beta=50 \text{ GB/s}$ .



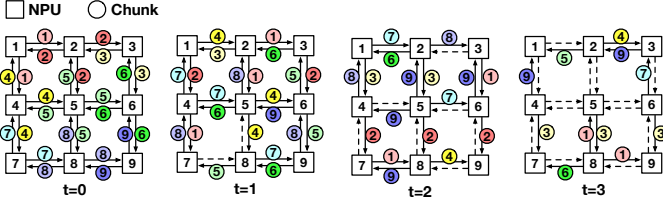


Fig. 14. An example All-Gather collective algorithm synthesized by TACOS over a homogeneous  $3 \times 3$  2D Mesh topology.

the send and receive operations at the link granularity. Each link is equipped with message queues and can process (i.e., send and receive) only one message at a time. For example, if two messages contend for the same link, only one message is sent out in a first-come, first-served order. Such simulation infrastructure enables the first-order modeling of network queuing and congestion effects for networks at scale.

**Validation.** We validated ASTRA-sim’s baseline All-Reduce implementations over two real systems: an NVLink-based, 8-GPU V100 server and a 32-TPUv3 ( $8 \times 4$ ) cluster. We observed ASTRA-sim to be 6–8% off on average, with the difference close to 1% for large message sizes [59].

## VI. EVALUATIONS

### A. TACOS Synthesis Result

We begin our evaluation with a visual representation of a TACOS-synthesized algorithm. Specifically, we synthesized an All-Gather algorithm for a  $3 \times 3$  2D Mesh network, as shown in Fig. 14. The synthesized algorithm effectively avoids network contentions, demonstrating TACOS’ ability to autonomously alleviate link congestion during synthesis.

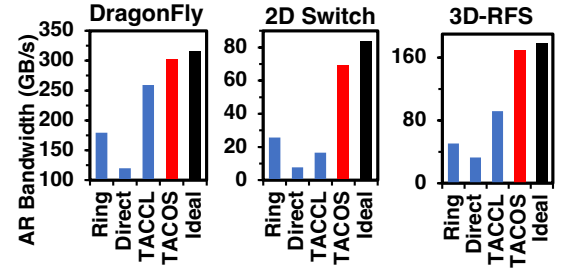
### B. Analysis of TACOS-Synthesized Algorithms

Modern distributed ML clusters typically employ symmetric topologies, whether homogeneous [2], [44], [45] or heterogeneous [1], [7], [18]. This section highlights the benefits of TACOS in handling these regular, symmetric topologies, compared to baseline collective algorithms and existing solutions.

**1. Topology Exploration.** We synthesized topology-aware All-Reduce algorithms for three distinct systems: a DragonFly topology ( $4 \times 5$ ) with link bandwidths of [400, 200] GB/s (per dimension), a 2D Switch ( $8 \times 4$ ) with [300, 25] GB/s, and 3D-RFS ( $2 \times 4 \times 8$ ) with [200, 100, 50] GB/s. Notably, the DragonFly topology is both asymmetric and heterogeneous, while Switch and 3D-RFS are symmetric topologies.

Fig. 15(a) illustrates the All-Reduce bandwidth (i.e., collective size  $\div$  collective time) synthesized by TACOS, alongside Ring, Direct, and TACCL-synthesized algorithms. TACOS effectively synthesizes collective algorithms, resulting in an average speedup of  $2.56\times$  over the baseline algorithms, including symmetric topologies. Moreover, thanks to its congestion-free search mechanism, TACOS even outperforms TACCL, consistently achieving more than 90% efficiency compared to the theoretical upper bound.

The heat map in Fig. 15(b) shows the average link utilization across all links for the DragonFly and 3D-RFS networks.



(a) All-Reduce bandwidth on heterogeneous topologies

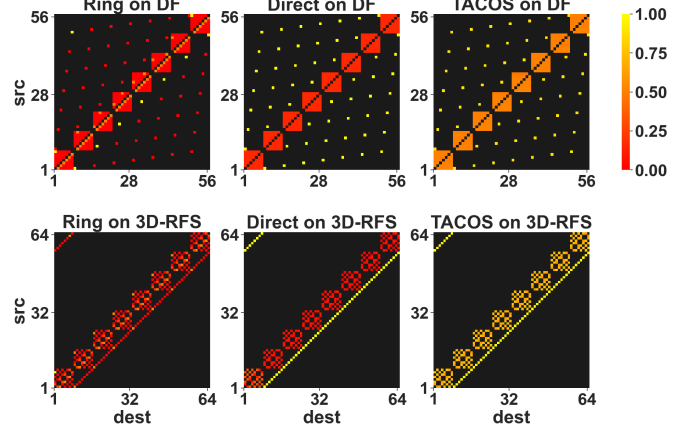


Fig. 15. (a) All-Reduce bandwidth of Ring and Direct basic algorithms, and TACCL and TACOS-synthesized algorithms. (b) Average link utilization measured on DragonFly and 3D-RFS. TACOS achieved 90.84% efficiency compared to the theoretical ideal.

TABLE V  
ALL-REDUCE COLLECTIVE TIME (WITH SYNTHESIS TIME IN PARENTHESES FOR BOTH TACOS AND TACCL) FOR A MULTI-NODE 3D-RFS SYSTEM WITH 2 TO 16 NODES (16 TO 128 NPUS), NORMALIZED OVER TACOS. ON AVERAGE, TACOS ACHIEVED 75.88% EFFICIENCY COMPARED TO THE THEORETICAL IDEAL.

#NPUs (#Nodes)	TACOS (ms)	TACCL (ms)	Ring	RHD	Direct	Ideal
16 (2)	1 (0.63)	2.89 (790)	7.14	5.27	4.04	1.00
32 (4)	1 (8.90)	4.10 (2001)	5.10	4.42	7.86	0.72
64 (8)	1 (97.92)	4.27 (7016)	4.80	5.83	16.84	0.68
128 (16)	1 (1080)	-	4.82	9.85	36.02	0.68

Topology-unaware basic algorithms lead to significant over-subscription on certain links, leaving others underutilized. However, the topology-aware collective algorithms synthesized by TACOS improve overall link utilization, distributing traffic more evenly across all links.

**2. Multi-node Analysis.** To demonstrate TACOS’ applicability for multi-node, symmetric AI clusters, we focus on the 3D-RFS topology and scale it by increasing the size of the last dimension (i.e., node). For instance, a 16-node system has  $2 \times 4 \times 16 = 128$  NPUs. Table V provides a comprehensive summary of these evaluations. For a 128-NPU cluster, TACCL incurred intractable synthesis times due to its NP-hard approach. On average, the collective algorithm generated by

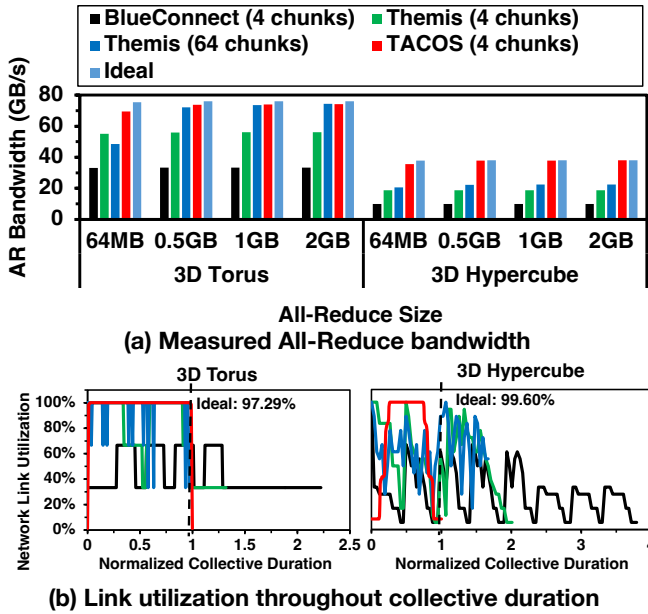


Fig. 16. (a) All-Reduce bandwidth comparison of TACOS against BlueConnect and Themis. TACOS achieves an average of 97.00% efficiency relative to the theoretical bound, regardless of collective size, while optimizations like Themis may struggle with latency-critical collectives. (b) Link utilization rate during collective execution, normalized by TACOS collective time. TACOS achieves 98.44% efficiency on average, regardless of topology shape.

TACOS exhibited a  $5.39\times$  speedup over the Ring baselines across all configurations. Direct baselines suffered as the Direct algorithm caused heavy network contention across the topology. TACOS achieved 75.88% efficiency relative to the theoretical ideal.

**3. BlueConnect and Themis.** To evaluate TACOS' performance against manually designed topology-aware collectives, we compare its results with the BlueConnect [57] and Themis [18] algorithms. The BlueConnect algorithm is designed for executing All-Reduce on symmetric multi-dimensional networks by applying Reduce-Scatter and All-Gather sequentially across each network dimension. Themis improves BlueConnect's efficiency by allowing distinct chunks to traverse network dimensions in an arbitrary manner, thereby balancing the load. We assessed TACOS' performance over a 3D Torus topology with  $\alpha=0.7 \mu s$  and  $1/\beta=25$  GB/s links.

We measured TACOS' benefits over Themis' target topology, a symmetric 3D Torus. The All-Reduce bandwidth is summarized in Fig. 16(a). TACOS achieved 95.90% efficiency relative to the ideal case. Themis, using 64 chunks, achieved similar efficiency for large collectives, but this came at the cost of latency. For latency-critical small collectives, Themis' efficiency dropped to 64.37%. When only 4 chunks per collective were used to avoid this issue, TACOS consistently outperformed Themis. Furthermore, for Themis-unfriendly asymmetric topologies like the 3D Hypercube, Themis achieved only 49.43% efficiency, while TACOS reached 98.10% of the ideal case. This discrepancy arises because Themis cannot alter the path each chunk takes, limiting its ability to mitigate network

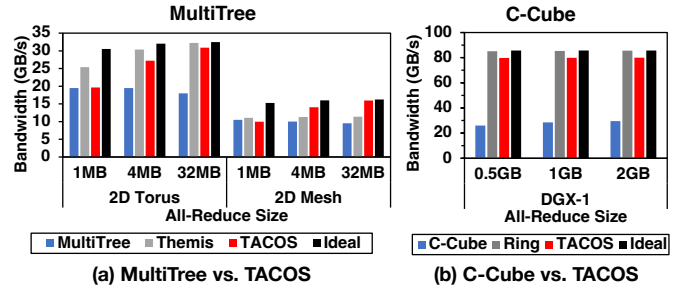


Fig. 17. All-Reduce bandwidth of TACOS and relative speedups compared to (a) MultiTree-synthesized and (b) C-Cube algorithms. TACOS consistently achieves near-optimal synthesis results, particularly for bandwidth-bound large collectives, resulting in an overall average efficiency of 85.25% relative to the ideal case.

contention when the optimal collective path for each network dimension is unknown. Consequently, TACOS demonstrated an average All-Reduce bandwidth  $2.01\times$  higher than Themis. Fig. 16(b) illustrates the link utilization during the collective execution. For the symmetric 3D Torus, both TACOS and Themis achieved nearly 100% utilization. However, in the 3D Hypercube topology, Themis experienced significant fluctuations in utilization due to network contention, which TACOS successfully avoided. We emphasize that TACOS achieved these results autonomously, without requiring any manual design efforts.

**4. MultiTree.** We conducted a comparison between TACOS and MultiTree, a spanning tree construction-based collective synthesizer [29]. The evaluation used 2D Torus and 2D Mesh topologies with  $\alpha=0.15 \mu s$  and  $1/\beta=16$  GB/s links. The summarized results are presented in Fig. 17(a), alongside corresponding Themis results and the ideal upper bound. Initially, TACOS demonstrated comparable performance for 1 MB All-Reduce, but outperformed as the collective size increased, exhibiting an average speedup of  $1.32\times$ . This is because, unlike TACOS, which can schedule and overlap multiple chunks concurrently for maximized network resource utilization, MultiTree does not support chunk-level overlap [26]. For larger collectives with multiple chunks, MultiTree cannot exploit full network bandwidth, leading to diminished results. Consequently, MultiTree saturated after 1 MB All-Reduce, while Themis and TACOS continued overlapping chunks, achieving 89.93% and 92.15% of theoretical efficiencies, respectively. Since Themis does not yield optimal results for asymmetric networks, for 2D Mesh, TACOS demonstrated a notable 82.60% efficiency. This underscores TACOS' ability to synthesize near-optimal collectives autonomously, regardless of target topology or collective size.

**5. C-Cube.** C-Cube [27] manually lays out two binary trees over a DGX-1 topology and schedules two tree-based collectives concurrently. We modeled an equivalent DGX-1 topology with  $\alpha=0.7 \mu s$  and  $1/\beta=25$  GB/s links. The comparison is illustrated in Fig. 17(b), where TACOS achieved  $2.86\times$  better performance on average. This can be attributed to C-Cube disabling 2 out of 6 available links per NPU to map two contention-free tree routes (as shown in Fig. 10(c) of [27]).

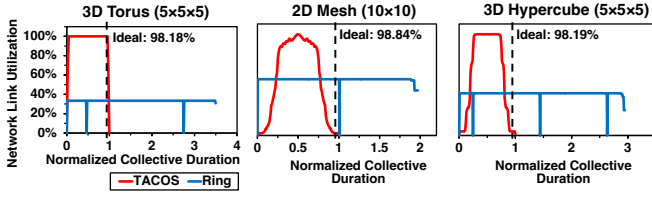


Fig. 18. Network utilization of TACOS-synthesized and Ring algorithms during the execution of an All-Reduce. The collective duration is normalized by TACOS collective time. On average, TACOS achieved 98.40% efficiency compared to the theoretical ideal.

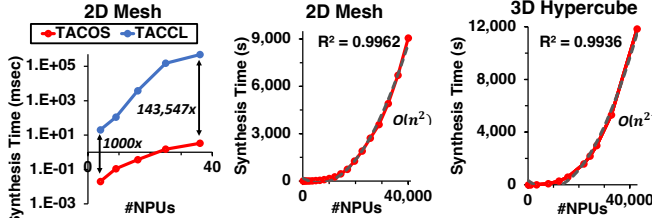


Fig. 19. Synthesis time of TACOS and TACCL for various-sized homogeneous 2D Mesh and 3D Hypercube topologies, using TACOS with 64 parallel threads. TACOS synthesis time exhibited  $O(n^2)$  complexity, where  $n$  is the number of NPUs.

Additionally, the remaining 4 links are not always fully utilized due to the tree-based approach for All-Reduce (i.e., links remain idle when there are no chunks to overlap), further reducing effective network utilization. Compared to the ideal bound, C-Cube achieved 32.63% efficiency, while TACOS reached 93.26%, comparable to the Ring baseline at 99.61%.

**6. Asymmetric Topologies.** To illustrate TACOS' efficacy over asymmetric networks, as shown in Fig. 18, we measured link utilization during the execution of an All-Reduce collective over the homogeneous 3D Torus ( $5 \times 5 \times 5$ ), 2D Mesh ( $10 \times 10$ ), and 3D Hypercube ( $5 \times 5 \times 5$ ) topologies. While 2D Mesh and 3D Hypercube are asymmetric, 3D Torus is symmetric and shown for comparison. We also included the utilization results of the Ring algorithm. It is noteworthy that the symmetric 3D Torus is highly suitable for collective communications, with TACOS achieving 100% utilization throughout the execution. However, the inherent asymmetry of 2D Mesh and 3D Hypercube topologies introduces inefficiencies, as not all chunks can depart or arrive simultaneously, causing some links to be idle at the start or end of execution. This effect is visible at  $t=3$  in Fig. 14. Nonetheless, TACOS successfully maximizes link utilization after saturation. Across all scenarios, TACOS achieved 98.40% efficiency relative to the theoretical ideal.

### C. Scalability Analysis

To demonstrate TACOS' scalability, we synthesized All-Reduce algorithms for homogeneous 2D Mesh and 3D Hypercube topologies with up to 43K NPUs, utilizing 64 parallel threads, and measured the synthesis time. The results are summarized in Fig. 19. Compared to TACCL for up to 36 NPUs, TACOS demonstrated significantly better scalability. Due to the NP-hard nature of ILP-based optimizers, the synthesis time

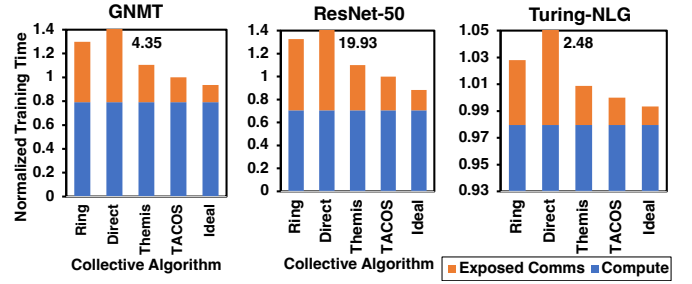


Fig. 20. End-to-end training time of GNMT, ResNet-50, and Turing-NLG. GNMT was trained on a 64-NPU 3D-RFS system. ResNet-50 and Turing-NLG are trained on a 3D-RFS with 32 nodes. All results are normalized over the corresponding TACOS result. TACOS achieved 82.38% communication efficiency and 93.61% end training efficiency compared to theoretical bounds.

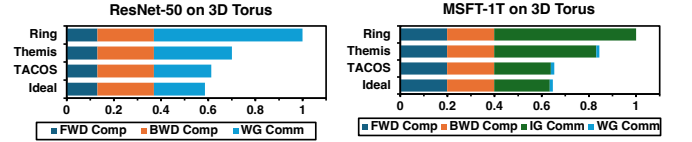


Fig. 21. End-to-end training time breakdown of ResNet-50 and MSFT-1T over a 1.024-NPU, 3D Torus. Forward and backward pass computation, and exposed input/weight gradient communications are shown. All results are normalized over the Ring. TACOS achieved 97.32% efficiency over the theoretical ideal scenario.

gap between TACOS and TACCL increased from  $10^3$  to  $10^5$  as the topology size increased by  $9 \times$ .

Further scalability analysis shows that TACOS synthesized an All-Reduce algorithm for a 2D Mesh topology with 40K NPUs in 2.52 hours and for a 3D Hypercube topology with 43K NPUs in 3.29 hours. The synthesis time scales quadratically with the number of NPUs, i.e.,  $O(n^2)$ , where  $n$  is the number of NPUs. This indicates that TACOS' synthesis time is linear relative to the search space, which consists of  $O(n)$  chunks and  $\Theta(n)$  links.

### D. End-to-End Application

Lastly, to elucidate TACOS' implications in distributed ML, we conducted an evaluation of the end-to-end training performance of GNMT [60], ResNet-50 [61], and Turing-NLG [62], representative models of vision and LLM workloads. For models employing data parallelism, communication becomes exposed at the end of each training iteration [18]. GNMT was evaluated on a small, 8-node (64 NPU) 3D-RFS topology, whereas ResNet-50 and Turing-NLG were run over a larger 32-node (256 NPU) cluster. The normalized training time over TACOS is illustrated in Fig. 20. Also, to model and showcase TACOS' applicability for larger ML clusters, we evaluated ResNet-50 and MSFT-1T [6] over a symmetric and homogeneous 3D Torus network with 1,024 NPUs. The breakdown of the normalized training time is depicted in Fig. 21.

Across all workloads, leveraging TACOS resulted in enhancements of  $1.58 \times$  and  $1.21 \times$  in end-to-end training performance over the baseline Ring and optimized Themis algorithms, respectively. Notably, compared to the theoretical upper bound, TACOS achieved the communication efficiency of

93.17%, thereby yielding an end-to-end efficiency of 97.32% compared to the ideal case.

## VII. RELATED WORK

### A. Time-expanded Networks

TEN is being leveraged in various fields that require communication or flow-based optimization. Notable examples include vehicle traffic management [54], logistics [39], and traffic signal optimization [55]. However, to the best of our knowledge, this work is the first attempt to bring TEN to distributed ML and collective communication with adequate optimization approaches.

### B. Collective Communications

**Collective Algorithm Representation.** Some works have proposed domain-specific languages for depicting collective algorithms [63], [64]; however, these are oriented toward the manual design of human-crafted collective algorithms. Other works [19], [20] use custom representations.

**Co-optimization of Collective and Topology.** Previous efforts have addressed topology optimization for specific collective algorithms [65], [66]. Additionally, the co-optimization of collectives and network topology remains an active research domain [29], [67]. Given TACOS’s ability to perform synthesis on arbitrary topologies, it can harmonize with such endeavors.

**Collective Algorithm Manual Design.** Designing specialized collective algorithms for specific topologies is an ongoing topic of interest. Recent examples include TTO for 2D Mesh [26] and PAARD for DragonFly [49]. However, this approach requires engineering and validation efforts for each topology variant, which can be automated or augmented with synthesizers like TACOS.

### C. Collective Algorithm Synthesizers

**Solver-based.** SCCL [20] synthesizes latency- and bandwidth-optimal collective algorithms by leveraging a satisfiability solver. To achieve this, SCCL captures the design space of collective communication in linear equations. To derive the linear relationships, SCCL assumes a k-synchronous collective algorithm, i.e., all NPUs executing identical transmissions clearly separated into steps and rounds. While SCCL guarantees optimality, this assumption only holds for homogeneous, symmetric, single-node networks, meaning the SCCL approach cannot be extended beyond these scenarios. TACCL [19] overcomes SCCL’s limitations through an ILP approach instead of a satisfiability problem, removing the k-synchronous assumption. However, the ILP approach requires TACCL to capture the entire search space in a number of linear equations. Since network congestion effects cannot be captured in linear equations, they are completely ignored in the formulation. Also, because the equations must be predefined, TACCL assumed a 2D network and constructed formulations solely based on it to model heterogeneity. Although this assumption can be extended to multi-dimensional topologies, the number of cross-dimensional constraints would increase, further complicating the ILP search space. Therefore, while

ILP guarantees the optimal solution of the provided formulation, as the equations themselves are limited, TACCL does not guarantee the optimality of the synthesized collective performance. Furthermore, TACCL had to assume the network is symmetrical to reduce the search space. With all these assumptions, both SCCL and TACCL synthesis efforts still scaled to only tens of NPUs.

**Tree-based.** Blink [28] constructs one-to-many spanning trees from a root to execute reduction and broadcast operations, effectively performing All-Reduce. The spanning tree construction process only takes network connectivity into account, thereby yielding sub-optimal results for heterogeneous topologies. To traverse the spanning tree upward and downward, all links are assumed to be bidirectional. Since the spanning tree is one-to-many, the approach is naturally detrimental for many-to-many collectives such as Reduce-Scatter or All-Gather, which are required by parallelization strategies like FSDP [41] or ZeRO [6]. TTO [26] optimizes Blink specifically for 2D Mesh topologies by constructing three spanning trees, all originating from edge NPUs. Although ideal for All-Reduce, TTO still shares the inefficiency for many-to-many collectives. MultiTree [29] generates height-balanced multiple spanning trees originating from all NPUs, making it suitable for many-to-many collectives. However, MultiTree only takes network connectivity into account, disregarding network heterogeneity. Unlike TACOS’s link-chunk matching, which automatically considers multiple chunk overlaps, MultiTree does not allow concurrent chunks to overlap. Therefore, while both may show similar synthesis results for small, latency-critical collectives, TACOS can achieve higher network utilization for larger collectives consisting of multiple chunks.

## VIII. CONCLUSION

In this paper, we underscore the importance of topology-aware collective algorithms in distributed ML and emphasize the necessity of having an autonomous collective algorithm synthesizer that supports arbitrary topologies. We introduce TACOS, an automated framework for orchestrating collective algorithm synthesis. TACOS supports diverse networks and demonstrates near-optimal link utilization while showcasing polynomial-time scalability.

## ACKNOWLEDGMENT

This work was supported through awards from Intel. Additionally, this research is supported by the Semiconductor Research Corporation (SRC), the SRC AIHW program, and the ACE Center for Evolvable Computing, one of the seven centers of the SRC JUMP 2.0 program. We also extend our sincere appreciation to Jiayi Huang and Le Qin for their assistance in helping us understand the details of MultiTree. We thank Ajaya Durg and Samvit Kaul for their constructive suggestions in implementing and evaluating this work. Finally, we greatly appreciate all the anonymous reviewers of this paper for dedicating their time and providing insightful comments to improve the quality of this work.



## APPENDIX

### ARTIFACT APPENDIX

#### A. Abstract

We open-source and provide TACOS, a topology-aware collective algorithm synthesizer, as the artifact proposed in this paper. TACOS is a C++17-based standalone program built upon a randomized link-chunk match-making algorithm. Therefore, TACOS supports any execution environment that can compile and execute C++17 code. TACOS uses CMake as its build automation tool, which is its only software dependency. Since the artifact is based on a randomized algorithm, the execution results may vary, but the trend of the target metric (e.g., collective algorithm bandwidth) should remain consistent.

#### B. Artifact Check-list (Meta-information)

- **Language:** C++17
- **Environment:** Any environment that can compile C++17
- **Metrics:** Collective Algorithm Bandwidth (GB/s)
- **Output:** Synthesized Collective Communication Algorithm
- **Disk space required:** A few MBs
- **Time to prepare workflow:** A few minutes
- **Time to complete experiments:** Depends on the search space size (target network and collective), but only a few seconds for  $O(100s)$  NPU.
- **Publicly available?** Yes
- **Code licenses:** MIT License
- **Archived:** DOI:10.5281/zenodo.13325902

#### C. Description

- 1) *How to Access:* <https://github.com/astra-sim/tacos>
- 2) *Hardware Dependencies:* None.
- 3) *Software Dependencies:* C++17 compiler and CMake.

#### D. Installation

The installation process is demonstrated in the README.md file included in the artifact (<https://github.com/astra-sim/tacos/blob/main/README.md>).

##### 1. Clone the TACOS repository:

- `git clone --recurse-submodules git@github.com:astra-sim/tacos.git`
- `cd tacos`

##### 2. Compile and run TACOS:

- `./tacos.sh`

## REFERENCES

- [1] Intel, “Gaudi Training Platform White Paper,” <https://habana.ai/wp-content/uploads/2019/06/Habana-Gaudi-Training-Platform-whitepaper.pdf>, 2019.
- [2] N. P. Jouppi, D. H. Yoon, G. Kurian, S. Li, N. Patil, J. Laudon, C. Young, and D. Patterson, “A Domain-Specific Supercomputer for Training Deep Neural Networks,” *Commun. ACM*, vol. 63, no. 7, p. 67–78, 2020.
- [3] NVIDIA, “NVIDIA H100 Tensor Core GPU,” <https://www.nvidia.com/en-us/data-center/h100>, 2022.
- [4] Cerebras Systems, “Cerebras Systems: Achieving Industry Best AI Performance Through A Systems Approach,” <https://cerebras.net/wp-content/uploads/2021/04/Cerebras-CS-2-Whitepaper.pdf>, 2021.
- [5] T. P. Morgan, “Inside Tesla’s Innovative And Homegrown Dojo AI Supercomputer,” <https://www.nextplatform.com/2022/08/23/inside-teslas-innovative-and-homegrown-dojo-ai-supercomputer>, 2022.
- [6] S. Rajbhandari, J. Rasley, O. Ruwase, and Y. He, “ZeRO: memory optimizations toward training trillion parameter models,” in *Proceedings of the 2020 International Conference for High Performance Computing, Networking, Storage and Analysis (SC ’20)*, 2020.
- [7] NVIDIA, “NVIDIA DGX SuperPOD: Instant Infrastructure for AI Leadership,” <https://resources.nvidia.com/en-us-auto-datacenter/nvpod-superpod-wp-09>, 2020.
- [8] J. Kim, W. J. Dally, S. Scott, and D. Abts, “Technology-Driven, Highly-Scalable Dragonfly Topology,” in *Proceedings of the 2008 International Symposium on Computer Architecture (ISCA ’08)*, 2008, pp. 77–88.
- [9] Groq, “GroqRack Compute Cluster,” <https://wow.groq.com/groqrack-compute-cluster>, 2024.
- [10] I. Cutress, “Analyzing Intel’s Discrete Xe-HPC Graphics Disclosure: Ponte Vecchio, Rambo Cache, and Gelato,” <https://www.anandtech.com/show/15188/analyzing-intels-discrete-xe-hpc-graphics-disclosure-ponte-vecchio>, 2019.
- [11] NVIDIA, “NVIDIA NVLink High-Speed GPU Interconnect,” <https://www.nvidia.com/en-us/design-visualization/nvlink-bridges>, 2022.
- [12] —, “InfiniBand Networking Solutions,” <https://www.nvidia.com/en-us/networking/products/infiniband>, 2024.
- [13] Y. Li, I.-J. Liu, Y. Yuan, D. Chen, A. Schwing, and J. Huang, “Accelerating Distributed Reinforcement Learning with In-Switch Computing,” in *Proceedings of the 46th International Symposium on Computer Architecture (ISCA ’19)*, 2019, p. 279–291.
- [14] A. Sapio, M. Canini, C. Ho, J. Nelson, P. Kalnis, C. Kim, A. Krishnamurthy, M. Moshref, D. R. K. Ports, and P. Richtárik, “Scaling Distributed Machine Learning with In-Network Aggregation,” in *arXiv:1903.06701 [cs.DC]*, 2019.
- [15] H. Mikami, H. Suganuma, P. Uchupala, Y. Tanaka, and Y. Kageyama, “Massively Distributed SGD: ImageNet/ResNet-50 Training in a Flash,” in *arXiv:1811.05233 [cs.LG]*, 2019.
- [16] A. Sapio, M. Canini, C.-Y. Ho, J. Nelson, P. Kalnis, C. Kim, A. Krishnamurthy, M. Moshref, D. Ports, and P. Richtárik, “Scaling Distributed Machine Learning with In-Network Aggregation,” in *Proceedings of the 18th USENIX Symposium on Networked Systems Design and Implementation (NSDI ’21)*, 2021, pp. 785–808.
- [17] Y. Li, I.-J. Liu, Y. Yuan, D. Chen, A. Schwing, and J. Huang, “Accelerating Distributed Reinforcement learning with In-Switch Computing,” in *Proceedings of the 46th Annual International Symposium on Computer Architecture (ISCA ’19)*, 2019, pp. 279–291.
- [18] S. Rashidi, W. Won, S. Srinivasan, S. Sridharan, and T. Krishna, “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 (ISCA ’22)*, 2022, p. 581–596.
- [19] A. Shah, V. Chidambaram, M. Cowan, S. Maleki, M. Musuvathi, T. Mytkowicz, J. Nelson, O. Saarikivi, and R. Singh, “TACCL: Guiding Collective Algorithm Synthesis using Communication Sketches,” in *Proceedings of the 20th USENIX Symposium on Networked Systems Design and Implementation (NSDI ’23)*, 2023, pp. 593–612.
- [20] A. Jangda, J. Huang, G. Liu, A. H. N. Sabet, S. Maleki, Y. Miao, M. Musuvathi, T. Mytkowicz, and O. Saarikivi, “Breaking the Computation and Communication Abstraction Barrier in Distributed Machine Learning Workloads,” in *Proceedings of the 27th ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS ’22)*, 2022, p. 402–416.
- [21] E. Chan, R. van de Geijn, W. Gropp, and R. Thakur, “Collective Communication on Architectures That Support Simultaneous Communication over Multiple Links,” in *Proceedings of the Eleventh ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP ’06)*, 2006, p. 2–11.
- [22] S. Rashidi, S. Sridharan, S. Srinivasan, and T. Krishna, “ASTRA-SIM: Enabling SW/HW Co-Design Exploration for Distributed DL Training Platforms,” in *Proceedings of the 2020 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS ’20)*, 2020, pp. 81–92.
- [23] R. Thakur, R. Rabenseifner, and W. Gropp, “Optimization of Collective Communication Operations in MPICH,” *Int. J. High Perform. Comput. Appl.*, vol. 19, no. 1, p. 49–66, 2005.
- [24] S. Jeageau, “Massively Scale Your Deep Learning Training with NCCL 2.4,” <https://developer.nvidia.com/blog/massively-scale-deep-learning-training-nccl-2-4>, 2019.

- [25] M. Cho, U. Finkler, M. Serrano, D. Kung, and H. Hunter, "BlueConnect: Decomposing all-reduce for deep learning on heterogeneous network hierarchy," *IBM Journal of Research and Development*, vol. 63, no. 6, 2019.
- [26] S. Laskar, P. Majhi, S. Kim, F. Mahmud, A. Muzahid, and E. J. Kim, "Enhancing Collective Communication in MCM Accelerators for Deep Learning Training," in *Proceedings of the 2024 IEEE International Symposium on High-Performance Computer Architecture (HPCA '24)*, 2024.
- [27] S. Cho, H. Son, and J. Kim, "Logical/Physical Topology-Aware Collective Communication in Deep Learning Training," in *Proceedings of the 2023 IEEE International Symposium on High-Performance Computer Architecture (HPCA '23)*, 2023, pp. 56–68.
- [28] G. Wang, S. Venkataraman, A. Phanishayee, N. Devanur, J. Thelin, and I. Stoica, "Blink: Fast and Generic Collectives for Distributed ML," in *Proceedings of the 2020 Machine Learning and Systems (MLSys '20)*, vol. 2, 2020, pp. 172–186.
- [29] J. Huang, P. Majumder, S. Kim, A. Muzahid, K. H. Yum, and E. J. Kim, "Communication Algorithm-Architecture Co-Design for Distributed Deep Learning," in *Proceedings of the 48th Annual International Symposium on Computer Architecture (ISCA '21)*, 2021, pp. 181–194.
- [30] T. Khan, S. Rashidi, S. Sridharan, P. Shurpali, A. Akella, and T. Krishna, "Impact of RoCE Congestion Control Policies on Distributed Training of DNNs," in *Proceedings of the 2022 IEEE Symposium on High-Performance Interconnects (HOTI '22)*, 2022, pp. 39–48.
- [31] R. Majumder and J. Wang, "DeepSpeed: Extreme-scale model training for everyone," <https://www.microsoft.com/en-us/research/blog/deepspeed-extreme-scale-model-training-for-everyone>, 2020.
- [32] NVIDIA, "NVIDIA Collective Communication Library (NCCL)," <https://developer.nvidia.com/nccl>, 2017.
- [33] AMD, "ROCCL 2.18.3 Documentation," <https://rocm.docs.amd.com/projects/rocl/en/latest>, 2023.
- [34] Microsoft, "Microsoft Collective Communication Library," <https://github.com/microsoft/msccl>, 2023.
- [35] Intel, "Intel oneAPI Collective Communications Library," <https://oneapi-src.github.io/oneCCL>, 2021.
- [36] E. Gabrielyan and R. Hersch, "Network topology aware scheduling of collective communications," in *Proceedings of the 10th International Conference on Telecommunications (ICT '03)*, 2003, pp. 1051–1058.
- [37] A. Paulus, M. Rolínek, V. Musil, B. Amos, and G. Martius, "CombOpt-Net: Fit the Right NP-Hard Problem by Learning Integer Programming Constraints," in *Proceedings of the 38th International Conference on Machine Learning (ICML '21)*, vol. 139, 2021, pp. 8443–8453.
- [38] E. Köhler, K. Langkau, and M. Skutella, "Time-Expanded Graphs for Flow-Dependent Transit Times," in *Proceedings of the 2002 European Symposium on Algorithms (ESA '02)*, 2002, pp. 599–611.
- [39] S. Belieres, M. Hewitt, N. Jozefowicz, and F. Semet, "A time-expanded network reduction matheuristic for the logistics service network design problem," *Transportation Research Part E: Logistics and Transportation Review*, vol. 147, 2021.
- [40] A. Tafreshian, M. Abdolmaleki, N. Masoud, and H. Wang, "Proactive shuttle dispatching in large-scale dynamic dial-a-ride systems," *Transportation Research Part B: Methodological*, vol. 150, pp. 227–259, 2021.
- [41] M. Ott, S. Shleifer, M. Xu, P. Goyal, Q. Duval, and V. Caggiano, "Fully Sharded Data Parallel: faster AI training with fewer GPUs," <https://engineering.fb.com/2021/07/15/open-source/fsdp>, 2021.
- [42] B. Klenk, N. Jiang, G. Thorson, and L. Dennison, "An In-Network Architecture for Accelerating Shared-Memory Multiprocessor Collectives," in *Proceedings of the 47th Annual International Symposium on Computer Architecture (ISCA '20)*, 2020, pp. 996–1009.
- [43] K. Kandalla, H. Subramoni, A. Vishnu, and D. K. Panda, "Designing topology-aware collective communication algorithms for large scale InfiniBand clusters: Case studies with Scatter and Gather," in *Proceedings of the 2010 IEEE International Symposium on Parallel & Distributed Processing, Workshops and Phd Forum (IPDPSW '10)*, 2010.
- [44] N. Jouppi, G. Kurian, S. Li, P. Ma, R. Nagarajan, L. Nai, N. Patil, S. Subramanian, A. Swing, B. Towles, C. Young, X. Zhou, Z. Zhou, and D. A. Patterson, "TPU v4: An Optically Reconfigurable Supercomputer for Machine Learning with Hardware Support for Embeddings," in *Proceedings of the 50th Annual International Symposium on Computer Architecture (ISCA '23)*, 2023.
- [45] A. Vahdat and M. Lohmeyer, "Enabling next-generation AI workloads: Announcing TPU v5p and AI Hypercomputer," <https://cloud.google.com/blog/products/ai-machine-learning/introducing-cloud-tpu-v5p-and-ai-hypercomputer>, 2023.
- [46] E. Sigler and B. Chess, "Scaling Kubernetes to 7,500 nodes," <https://openai.com/research/scaling-kubernetes-to-7500-nodes>, 2021.
- [47] M. O. Kevin Lee, Adi Gangidi, "Building Meta's GenAI Infrastructure," <https://engineering.fb.com/2024/03/12/data-center-engineering/building-metas-genai-infrastructure/>, 2024.
- [48] S. Kumar and N. Jouppi, "Highly Available Data Parallel ML training on Mesh Networks," in *arXiv:2011.03605 [cs.LG]*, 2020.
- [49] J. Ma, D. Dong, C. Li, K. Wu, and L. Xiao, "PAARD: Proximity-Aware All-Reduce Communication for Dragonfly Networks," in *Proceedings of the 2021 IEEE International Conference on Parallel & Distributed Processing with Applications, Big Data & Cloud Computing, Sustainable Computing & Communications, Social Computing & Networking (ISPA/BDCloud/SocialCom/SustainCom '21)*, 2021, pp. 255–262.
- [50] J. Kim, J. Balfour, and W. Dally, "Flattened Butterfly Topology for On-Chip Networks," in *Proceedings of the 40th Annual Symposium on Microarchitecture (MICRO '07)*, 2007, pp. 172–182.
- [51] M. Flajslik, E. Borch, and M. A. Parker, "Megafly: A Topology for Exascale Systems," in *High Performance Computing*, 2018, pp. 289–310.
- [52] M. Besta and T. Hoefer, "Slim Fly: A Cost Effective Low-Diameter Network Topology," in *arXiv:1912.08968 [cs.NI]*, 2020.
- [53] Y. Ajima, S. Sumimoto, and T. Shimizu, "Tofu: A 6D Mesh/Torus Interconnect for Exascale Computers," *Computer*, vol. 42, no. 11, pp. 36–40, 2009.
- [54] M. Ferrati and L. Pallottino, "A time expanded network based algorithm for safe and efficient distributed multi-agent coordination," in *Proceedings of the 52nd IEEE Conference on Decision and Control (CDC '13)*, 2013, pp. 2805–2810.
- [55] E. Köhler and M. Strehler, "Combining Static and Dynamic Models for Traffic Signal Optimization Inherent Load-dependent Travel Times in a Cyclically Time-expanded Network Model," *Procedia - Social and Behavioral Sciences*, vol. 54, pp. 1125–1134, 2012.
- [56] R. W. Hockney, "The communication challenge for MPP: Intel Paragon and Meiko CS-2," *Parallel Computing*, vol. 20, no. 3, pp. 389–398, 1994.
- [57] W. Won, S. Rashidi, S. Srinivasan, and T. Krishna, "LIBRA: Enabling Workload-Aware Multi-Dimensional Network Topology Optimization for Distributed Training of Large AI Models," in *Proceedings of the 2024 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS '24)*, 2024, pp. 205–216.
- [58] W. Won, T. Heo, S. Rashidi, S. Sridharan, S. Srinivasan, and T. Krishna, "ASTRA-sim2.0: Modeling Hierarchical Networks and Disaggregated Systems for Large-model Training at Scale," in *Proceedings of the 2023 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS '23)*, 2023, pp. 283–294.
- [59] ASTRA-sim, "ASTRA-sim Validation," <https://astra-sim.github.io/astra-sim-docs/validation/validation.html>, 2024.
- [60] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey, J. Klingner, A. Shah, M. Johnson, X. Liu, L. Kaiser, S. Gouws, Y. Kato, T. Kudo, H. Kazawa, K. Stevens, G. Kurian, N. Patil, W. Wang, C. Young, J. Smith, J. Riesa, A. Rudnick, O. Vinyals, G. Corrado, M. Hughes, and J. Dean, "Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation," in *arXiv:1609.08144 [cs.CL]*, 2016.
- [61] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in *Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR '16)*, 2016, pp. 770–778.
- [62] C. Rosset, "Turing-NLG: A 17-billion-parameter language model by Microsoft," <https://www.microsoft.com/en-us/research/blog/turing-nlg-a-17-billion-parameter-language-model-by-microsoft>, 2020.
- [63] M. Cowan, S. Maleki, M. Musuvathi, O. Saarikivi, and Y. Xiong, "MSC-CLang: Microsoft Collective Communication Language," in *Proceedings of the 28th ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '23)*, 2023, p. 502–514.
- [64] A. Jangda, J. Huang, G. Liu, A. H. N. Sabet, S. Maleki, Y. Miao, M. Musuvathi, T. Mytkowicz, and O. Saarikivi, "Breaking the Computation and Communication Abstraction Barrier in Distributed Machine Learning Workloads," in *Proceedings of the 27th ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '22)*, 2022, p. 402–416.

- [65] T. T. Nguyen and M. Wahib, “An Allreduce Algorithm and Network Co-design for Large-Scale Training of Distributed Deep Learning,” in *Proceedings of the 21st International Symposium on Cluster, Cloud and Internet Computing (CCGrid '21)*, 2021, pp. 396–405.
- [66] L. Zhao, S. Pal, T. Chugh, W. Wang, J. Fantl, P. Basu, J. Khoury, and A. Krishnamurthy, “Efficient Direct-Connect Topologies for Collective Communications,” in *arXiv:2202.03356 [cs.NI]*, 2023.
- [67] W. Wang, M. Khazraee, Z. Zhong, M. Ghobadi, Z. Jia, D. Mudigere, Y. Zhang, and A. Kewitsch, “TopoOpt: Co-optimizing Network Topology and Parallelization Strategy for Distributed Training Jobs,” in *Proceedings of the 20th USENIX Symposium on Networked Systems Design and Implementation (NSDI '23)*, 2023, pp. 739–767.