# MAD-Max Beyond Single-Node: Enabling LargeMachine Learning Model Acceleration onDistributed Systems

Samuel Hsia<sup>1,2</sup>, Alicia Golden<sup>1,2</sup>, Bilge Acun<sup>1</sup>, Newsha Ardalani<sup>1</sup>, Zachary DeVito<sup>1</sup>, Gu-Yeon Wei<sup>2</sup>, David Brooks<sup>2</sup>, Carole-Jean Wu<sup>1</sup>

<sup>1</sup>FAIR at Meta, <sup>2</sup>Harvard University

shsia@g.harvard.edu, carolejeanwu@meta.com

Abstract—Training and deploying large-scale machine learning models is time-consuming, requires significant distributed computing infrastructures, and incurs high operational costs. Our analysis, grounded in real-world large model training on datacenter-scale infrastructures, reveals that  $14\sim32\%$  of all GPU hours are spent on communication with no overlapping computation. To minimize this outstanding communication latency and other inherent at-scale inefficiencies, we introduce an agile performance modeling framework, MAD-Max. This framework is designed to optimize parallelization strategies and facilitate hardware-software co-design opportunities. Through the application of MAD-Max to a suite of real-world large-scale ML models on state-of-the-art GPU clusters, we showcase potential throughput enhancements of up to  $2.24\times$  for pretraining and up to  $5.27\times$  for inference scenarios, respectively.

#### I. INTRODUCTION

Billion-parameter large language models (LLMs) [9], [49], [61], [62] power applications that have shown far-reaching impact across different domains [15], [16], [38], [48]. Similarly, trillion-parameter recommendation models [40], [72] have demonstrated state-of-the-art user modeling and content understanding across search [6], [11], [31], [76], social media [1], [18], [19], [71], e-commerce [78], [79], and entertainment [20]. As these large-scale ML models increase in size and complexity [18], [19], the corresponding training and inference workloads become ever more resource-intensive. Without efficient mappings between these large-scale ML workloads and their underlying distributed systems, model training and exploration can easily require millions of GPU hours, levying high operational costs, compute resource requirements, and energy consumption [9], [61], [62].

Figure 1 shows the projected resource-performance pareto frontier of training a state-of-the-art deep learning recommendation model (DLRM) using default workload-system mapping strategy on public cloud instances. In this case, we quantify compute resource requirements with aggregate GPU hours per 1 billion samples, where aggregate GPU hours of different generations of GPUs are normalized based on the A100's peak FLOPS. Further improving upon this

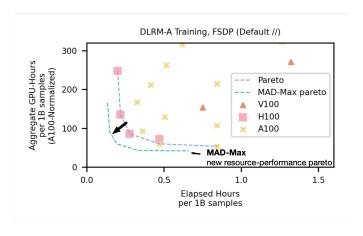


Fig. 1. Our performance model – MAD-Max – improves upon the resource-performance pareto frontier of large-scale ML workloads by identifying new hardware-software mappings and solutions.

resource-performance pareto frontier requires researchers to take into account underlying distributed systems [14], [28]–[30], [42], [43], [45], [46] and how we map models and tasks onto underlying distributed systems – parallelization strategy. In this paper, we propose a distributed ML performance model – *MAD-Max* – for identifying potential avenues for improvement (green, dotted line). Nonetheless, pinpointing the specific distributed systems and parallelization strategies needed for realizing these improvements in performance and operational compute resource requirements is challenging, as evidenced by the three general approaches for optimizing runtime performance of large ML models.

The first option involves applying industry-standard parallelization strategies (Figure 1: blue, dotted line) that target feasibility without fully optimizing hardware usage (e.g., FSDP [54], [75]). The second option is to custom-design custom hierarchical parallelization strategies specific to the model, task, and system [59]. This maximizes hardware efficiency but is complex from an engineering standpoint and not easily adaptable across different tasks. The third option is to use

software tools to predict system performance before training or deployment, though existing tools have several limitations, such as being training-specific or hardware architecture-dependent. To address the need for an agile exploration tool to identify parallelization strategies tailored to different use-cases, we introduce our distributed ML performance model – *MAD-Max* – and evaluate it on a suite of real-world, large ML models, including deep learning recommender systems and LLMs [8]–[10], [13], [41], [51], [61], [62], [76].

In this work, we first characterize a suite of real-world, large ML models at both model- and datacenter-deployment scales (Section III). At the model architecture level, we identify performance-critical hardware requirements based on the models' compute and memory characteristics. At the datacenter scale, we quantify the required communication by conducting a fleet-wide training characterization, revealing that  $14\sim32\%$  of all GPU hours are spent on *communication with no concurrent computation* (i.e., exposed communication).

To enable agile exploration of the parallelization design space, MAD-Max first estimates the system performance of large-scale ML workloads. The performance model takes in target ML model architecture, task details, parallelization scheme, and distributed system hardware to generated perdevice traces. These per-device traces are then pieced together to estimate the overall system performance of the target ML model and task. Additionally, the performance model generates detailed breakdowns of both communication collectives and computation-communication overlap efficiency, enabling users to identify future optimization opportunities. Our performance model is validated against multiple real-world large-scale distributed training experiments, demonstrating 97% and 91% performance prediction accuracies on serialized and overlapped execution, respectively.

Using MAD-Max, we identify parallelization strategies that result in throughput improvements across our suite of large ML models – achieving up to  $2.24\times$  and  $5.27\times$  throughput improvements for pre-training and inference, respectively. By extending our analysis to parallelization strategies that are not constrained by the memory capacities of existing training platforms, we discover strategies capable of delivering up to  $2.43 \times$  and  $12.13 \times$  throughput improvement for pre-training and inference, respectively. Furthermore, MAD-Max provides critical insights on how model-level compute and communication requirements alter optimal parallelization strategy and increasing LLM context lengths calls for solutions beyond purely parallelization exploration (Section VI). We also study how different generations of GPUs and other commodity hardware platforms impact overall training efficiency and follow up with a future technologies scaling study by showing the effects of asymmetrically improving systems components like compute efficiency, memory capacity and bandwidth, and hierarchical interconnect bandwidth (Section VI).

The main contributions of this work are as follows:

 We propose a performance model that enables agile exploration of the distributed ML training and deployment design space. Our performance model targets both

- implemented and future models alike, enabling accurate performance estimation with different model architectures, tasks, hardware devices, and distributed systems.
- We show model-level insights on how parallelization strategies interact with DLRM and its transformer and mixture-of-experts variants. We show how asymmetric compute and communication requirements from transformer and mixture-of-experts components lead to different optimal parallelization strategies. Additionally, we demonstrate the limits of solely optimizing parallelization strategies on LLMs of increasing context length.
- We show that to improve large ML model training and inference throughput, hardware specifications across compute, memory, and interconnect have to be concurrently improved.

We have open-sourced MAD-Max and sample experiment to enable follow-on work for modeling the interaction between parallelization strategies, models, tasks, and distributed systems on ML system performance.

# II. BACKGROUND

In this section, we introduce a suite of model architectures across both recommender systems and LLMs. We then outline three tasks for these models: pre-training, fine-tuning, and inference (Section II-A). Lastly, we discuss the parallelization strategies currently used to map the workloads (i.e., model and task) onto the distributed systems (Section II-B).

### A. Models and Tasks

Deep learning based recommender systems and LLMs follow the general model architecture of representing categorical inputs as embedding vectors and then processing these embedding vectors with model-specific computation layers. This means that there are many shared components that are emphasized to different degrees by each model: embedding tables, multilayer perceptrons (MLPs), and more intricate dense processing layers like transformer blocks. We focus on the following five classes of models throughout the paper:

- 1) **DLRM.** The canonical at-scale recommendation model takes in dense and sparse features. Dense features, such as, user age and current time, are processed by MLP layers while sparse categorical features are processed as lookups into large embedding tables. These results are then fed into a feature interaction layer, where these intermediate values are either concatenated or multiplied with one another via dot products [64], [65]. The result of this feature interaction layer is then fed into MLP layers to generate predictions like Click-Through Rate (CTR) [41]. For many large-scale DLRM models, storing and communicating trillion-parameter scale embedding tables is the primary system bottleneck [17], [18], [22], [23], [33], [34], [40], [58], [67].
- DLRM-Transformer. As sparse features for recommendation models increase in complexity, corresponding model architectures have also evolved to better model implicit relationships between sparse features.

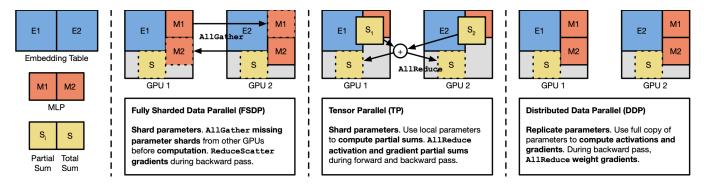


Fig. 2. For recommendation models, applying FSDP, TP, or DDP on an MLP layer requires either sharding or replicating parameters and communicating either parameters (orange) or partial sums (yellow). In this example, the embedding table's prohibitively large capacity requires it to be sharded.

Some DLRM variants replace concatenation and dotproduct based feature interactions with transformer encoder layers that model higher-order interactions and sequential relationship between sparse features. Others use transformer-style feature interaction layers to tackle challenges like behavior sequence modeling and personalized re-ranking [10], [51], [72]. From a systems perspective, transformer layers increase both compute and computation-communication overlap opportunities.

- 3) **DLRM-MoE.** In the context of DLRMs, applying Mixture-of-Experts (MoE) creates parallel Top MLPs that are conditionally activated based on feature interactions [76]. Because only a fraction of experts are active for each sample, DLRM-MoE increases model capacity and expert-to-expert communication while scaling computation at a lower rate.
- 4) **LLM.** Large language models (LLMs) also use the "look up embeddings then process them" architecture [9], [21], [53], [61], [74]. However, instead of using user and content categorical features, LLMs convert to-kens character sequences to input embeddings. Subsequent processing layers use alternating self-attention and feed-forward layers [63]. Unlike DLRMs, advancements in LLM modeling have been more focused on the processing layers than embeddings, reinforcing the importance of compute in LLM execution.
- 5) **LLM-MoE.** In the context of LLMs, one way to apply MoE is to replace the feed-forward layer in transformer blocks with experts. By applying this technique, the FLOPs per token will grow at a slower rate than overall model capacity while scaling up the model, leading to enabling efficient training and inference. While FLOPs becomes less of a concern, the non-blocking expert-to-expert communication that can be present during training presents systems challenges.

In terms of tasks, we are interested in pre-training, finetuning and inference. Pre-training stresses all of compute, memory capacity, and communication as it involves both forward and backward passes – along with retaining intermediate activations from the forward pass. The requirements of finetuning are a subset of pre-training, as the frozen parameters of a model do not require updates. Inference only requires the forward pass so compute is usually proportionally larger.

### B. Parallelization Strategies

A model layer can be either replicated or sharded across devices. We explore the following parallelization strategies (Figure 2 illustrates forward pass execution):

- 1) Fully Sharded Data Parallelism (FSDP). Parameters are *sharded* across devices. Before forward and backward pass, missing parameter shards are gathered via AllGather. During backward pass, weight gradients are reduced and sharded via ReduceScatter.
- 2) Tensor Parallelism (TP). Parameters are sharded across devices. During forward pass, each device uses its parameter shard to compute partial sums that are then aggregated via AllReduce. Same principle is applied for backward pass for gradients.
- 3) **Distributed Data Parallelism (DDP).** Parameters are *replicated* across devices. During forward pass, each device acts independently for computation. During backward pass, devices AllReduce weight gradients.

We apply one parallelization strategy for each layer type. Figure 2 depicts applying different parallelization strategies on an MLP layer and vanilla **model parallel** (**MP**) sharding for the embedding tables. Additionally, parallelization strategies can be applied hierarchically for multi-node systems, creating *N*-D parallelism strategies.

# III. CHARACTERIZATION

In this section, we first characterize a suite of real-world large ML models with respect to model capacity, parameter breakdowns, FLOPs, and memory bandwidth characteristics (Section III-A). To get a better understanding of the models' communication requirements, we conduct a fleet-wide characterization of at-scale training experiments (Section III-B).

# A. Individual Model Characterization

We first quantify the difference in compute, memory capacity, and bandwidth requirements between six real-world recommendation models and LLMs: DLRM-{A, B, C}, GPT-3 175B, LLaMA-65B, LLaMA 2-70B. Figure 3 quantifies this diversity of requirements with two key observations:

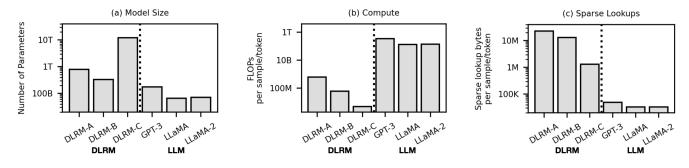


Fig. 3. For large ML models, the requirements for key system resources - (a) capacity, (b) compute, (c) bandwidth - vary by orders of magnitude.

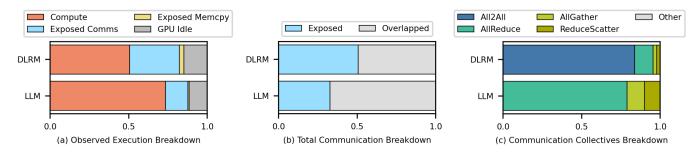


Fig. 4. (a) Compute and exposed communication make up the majority of observed at-scale training cycles. (b) The degree of communication overlapped with computation and data loading is workload dependent. Higher degree of overlap indicates better latency hiding of communication collectives. (c) Breakdown of communication collectives also varies by workload.

O1: Parameter count – and allocation across model layers – varies by orders of magnitude between models, impacting system capacity requirements. Recommendation models contain significantly more parameters than LLMs (Figure 3 (a)). Despite variation in parameter count across LLMs, GPT-3 consists of roughly 2–68× fewer parameters as compared to recommendation models. Training and deploying these recommendation models and LLMs require multi-node distributed systems, yet the size of the target model governs how many devices (i.e., GPUs) are required to fit the entire model and the viable set of scale-out parallelization strategies.

Additionally, virtually 100% of parameters in recommendation models are used for embeddings while almost all LLM parameters are dedicated to compute. This reflects the transformer-heavy computation of current LLM architectures, in contrast to embedding-driven, recommendation model architectures for at-scale personalization.

<u>O2</u>: Recommendation models require fewer FLOPs per sample as compared to LLMs, yet require >20× higher memory bandwidth for sparse lookups. Figures 3 (b, c) illustrate how recommendation models and LLMs show opposite trends for compute requirements and sparse lookup bandwidth. Sparse lookup bandwidth requirements for recommendation models far surpass LLMs – a fact that is consistent with how recommendation models have a higher proportion of parameters dedicated to embeddings. However, the opposite is true for compute requirements, as LLMs require significantly higher FLOPs per sample. As discussed in Section IV, these varying system requirements play an important role in the

design of an optimal parallelization strategies.

#### B. Fleet-wide Communication Characterization

In addition to model-level characterization, we look at fleetwide model training. We observe, over an extended period of time, the importance of communication for training the latest DLRM-style models and LLMs. Figure 4 quantifies the role of communication with two key observations:

O3: Compute and exposed communication make up the majority of observable training GPU cycles. Compute, defined as cycles with either device computation or memory lookups (orange) and exposed communication, defined as cycles with only inter-device communication (blue), make up >82% of all observable training GPU cycles for both DLRM and LLMs (Figure 4 (a)). The rest of the cycles are attributed to host-device communication – exposed memcpy (yellow) – and inactivity due to data ingestion, kernel launch overhead, etc. – GPU idle (grey). From this observation, we focus our performance modeling efforts on predicting the expected behavior of compute and communication cycles.

 $\underline{O4:}$  Differences in model architectures and parallelization strategies impact both the amount of compute-communication overlap and the types of communication collectives used. When model training spans multiple devices, replicating or sharding model components leads to communication calls involving parameters, activations and/or gradients. Being able to overlap these communication calls with computation so that the hardware devices are doing useful work is important for utilization. Figure 4 (b) shows that  $\sim 50\%$  of communication calls for DLRM training are overlapped

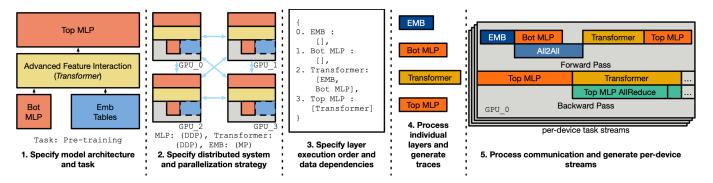


Fig. 5. Our performance model works in five stages. After workload specifications and layer execution orders are established, traces for individual layer execution are generated and then combined with required communication collectives to form complete computation and communication streams.

with computation, whereas >65% of communication calls for compute-dominated LLMs are overlapped.

Figure 4 (c) shows the spread of different communication collectives during training. For DLRM models, All2All is heavily emphasized while LLMs spend the majority of their communication cycles on AllReduce. This is a direct result of model architecture difference, and thus active parallelization strategy. Since DLRMs require large amounts of sparse lookups from sharded embedding tables, the per-device unique embedding lookups have to be distributed to each device via All2All. On the contrary, LLMs have fewer parameters and are more amenable to replication of compute parameters, allowing for DDP opportunities that require AllReduce for aggregating weight gradients.

In this section, we characterize real-world large ML models from model architecture and distributed training perspectives. From Section III-B we see that model architectures and the way in we map them onto distributed systems significantly impacts system resource utilization, and thus overall performance. To better understand how to best map current and future large ML models onto different distributed systems, we propose an agile, at-scale accurate performance model.

# IV. PROPOSED DESIGN

In this section, we outline the structure of our performance model – MAD-Max – for simulating distributed ML workloads. We begin with an overview of the model's design and the key assumptions it relies on, highlighting the role of execution traces in modeling the iterative behavior of large-scale ML tasks (Section IV-A). Then, we discuss how the model processes individual ML model layers according to their key characteristics (Section IV-B). We conclude by explaining the integration of these individually-processed layers into a unified computation and communication model that addresses the communication requirements dictated by the chosen parallelization strategy (Section IV-C).

# A. Design Overview

Our performance model, illustrated through a DLRM-Transformer case in Figure 5, is predicated on the notion that ML model layers, when treated as discrete blocks, can be used to create **per-device execution traces** for emulating the periteration behavior of distributed ML workloads. An "**execution trace**" in this context refers to a detailed record capturing the sequence and duration of both compute and communication events (i.e., streams) on each device. To simulate the periteration behavior of a distributed ML workload, MAD-Max constructs a dependency graph of layers, generates per-layer compute traces, and then pieces together the compute traces with traces of parallelization strategy specific communication collectives to form complete compute and communication streams. From per-iteration behavior, the performance model estimates overall throughput and other end-to-end serialized and overlapped execution breakdowns.

Users have to provide JSON files for: 1) model architecture via layer-specific configurations (e.g., number of MLP layers, embedding table dimension, number of transformer layers and heads), 2) distributed system specifications (e.g., Tensor Float (TF32) utilization, HBM peak bandwidth, AllReduce intranode interconnect utilization), and 3) task and parallelization strategy (e.g., pre-training/fine-tuning/inference, intra-/internode parallelization strategies).

With these configurations, individual layers are first processed by their primary system requirements. Examples include estimating embedding bag execution by the amount of embeddings to look up and per-GPU high-bandwidth memory (HBM) memory bandwidth and the time it takes to execute a transformer encoder layer by TF32 compute throughput. Based on the replication and sharding specified by the target parallelization strategy, the required communication collectives are processed by collective-specific intra- (e.g., NVLink) and inter- (e.g., Infiniband, RDMA over Converged Ethernet (RoCE)) node communication bandwidths.

We take into account task-level requirements (i.e., pretraining/fine-tuning/inference) to construct per-device computation and communication streams with data dependencies and potential computation-communication overlap.

#### **Assumptions:**

Since we focus on large-models, target distributed systems are multi-device in nature. For multi-device execution, a first-order analysis of execution behavior and overall performance can be estimated via modeling per-

node layer execution and inter-node parallelization communication. Kernel-level improvements (e.g., [47]), while not the focus of this work, can be effectively modeled as increased compute and memory lookup utilization.

- The performance model assumes that the entire model can be fit onto the training/inference devices (i.e., when sharded, the model can fit onto GPUs). Recent high-performance training platforms target this design point [40]. Design points where model parameters have to be shuffled back and forth between CPU and device are currently unsupported.
- Device-host communication (e.g., CPU-GPU data loading) is relatively a second-order consideration and mostly overlapped and hidden between training/inference iterations. This observation is shared in [40] and our fleet-wide characterization in Section III-B, Figure 4.

# B. Processing Individual Model Layers

Layers are processed by their main system requirement. For example, we illustrate how MLP and embedding bag performance are estimated differently in Figure 5.

**Compute Blocks.** Assuming that compute time is the main bottleneck for MLPs, we estimate compute time per layer as:

$$\sim$$
 (FLOPs per layer) / [(GPU peak FLOPS) \* Compute utilization]

where FLOPs per layer is determined by the MLP layer's dimensions and target batch size. GPU peak FLOPS are heavily dependent on data type (e.g., 32-bit, 16-bit FP/TF/BF) and whether or not tensor cores are enabled. We incorporate compute utilization – or in the case of GPUs, SM utilization/occupancy – as a factor in [0,1]. Typical compute utilization factors for A100s on layers in our models of interest are  $\sim$ 70%. We adopt a similar approach for modeling self-attention and fully-connected (FC) layers found in transformer layers, where FLOPs per layer is estimated by additional factors such as attention dimension and context length.

**Embedding Bags.** Assuming that lookup time is the main bottleneck for embedding bags, we estimate lookup time as:

where Lookup bytes is determined by the number of embedding tables, number of lookups per embedding table, embedding dimension, and embedding precision. Lookup bytes per GPU is highly parallelization strategy dependent. In this case, we assume that the embedding table is evenly sharded across GPUs in terms of both capacity and number of lookups. If the number of lookups are unevenly distributed between GPUs, we can adjust the lookup bytes per GPU on a per-GPU basis [58]. HBM utilization is a factor between [0,1] and typical values for embedding bags of interest are ~80% for A100s.

#### C. Piecing Together Computation and Comm. Streams

**Specifying Explicit Execution Order.** To generate perdevice traces for different ML tasks, an explicit execution priority must be established for the different layers. In Figure 5,

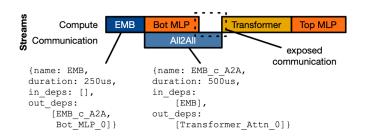


Fig. 6. Sample generated GPU compute and communication streams with labeled exposed communication.

we can establish the order as such (1) Embedding, (2) Bottom MLP, (3) Transformer, (4) Top MLP. During backward pass, the execution order will be reversed. If the target task is fine-tuning, we also specify frozen layers, reducing unnecessary computation and communication of certain weight gradients.

Generating Parallelization-Specific Streams. An explicit execution order by itself is not enough to construct accurate streams. A target parallelization strategy is required to specify the required communication collectives. Explicit data dependencies, along with parallelization strategy determine the blocking/non-blocking nature of the communication calls. In Figure 5, MLP and transformer layers are distributed via DDP while embedding tables are distributed via sharding.

Figure 6 illustrates generated forward pass streams from our DLRM-Transformer example. We see that the traces are slotted into a compute stream and communication stream. Each trace will have dependencies that come explicitly from execution annotations and implicitly from underlying parallelization strategies. For example, EMB has an explicit output dependency of Bot\_MLP\_0 and implicit output dependency of EMB\_c\_A2A from sharding the embedding table. EMB\_c\_A2A is blocking since Transformer\_Attn\_0 needs EMB\_c\_A2A's results.

**Estimating Communication Collective Execution.** We estimate All2All execution as:

~ ("SendCount" Bytes per GPU) / (Effective All2All BW)

where "SendCount" Bytes per GPU is the average number of bytes sent by each GPU to every other GPU. "SendCount" Bytes per GPU is dependent on not only "Lookup bytes per GPU" but also the sharding degree and number of devices. Since the All2All NCCL implementation is composed of individual point-to-point Send() and Recv() calls, it is bound by the slowest level of interconnect [69]. Thus, for baseline DGX systems, Effective All2All BW is set as that of either Infiniband or RoCE (i.e., whatever interconnect fabric is used to connect nodes of GPUs). For other cases, like an 8-GPU system, Effective All2All BW may be NVLink BW.

Likewise, we can generate a similar set of traces for the backward pass. Since the MLP and transformer layers are parallelized via DDP, we have non-blocking AllReduce communication calls during the backward pass. The AllReduce calls are for aggregating per-layer weight gradients and are thus non-blocking (i.e., they are not on the critical path for

	Evaluation Metric	Measured Result	Performance Model Result	Modeling Accuracy (%)
	Serialized Iteration Time (ms)	67.40 ms	65.30 ms	96.89%
DLRM-A	% Communication Exposed (%)	82.37%	75.46%	91.62%
	Throughput (MQPS)	1.2 MQPS [40]	1.21 MQPS	99.17%
DLRM-B	Tilloughput (WQI 3)	3.4 MQPS [40]	3.06 MQPS	90%
	GPU Hours for 306k steps	1,022,361	863,397	84.66%
LLaMA-70B	(2048 A100s)	Hrs	Hrs	
	Days to Train 1.4T Tokens	20.83 Days [61]	19.21 Days	92.27%

TABLE I
VALIDATION OF VARIOUS FIRST-ORDER EXECUTION METRICS.

backpropagation). We estimate the non-blocking AllReduce calls for weight gradient calls as:

#### ~ ("SendBuffer" Bytes / GPU) / (Effective AllReduce BW)

where "SendBuffer" Bytes is the total number of bytes sent by each GPU and is directly proportional to the number of parameters in each layer. Effective AllReduce BW is a ratio of intra-node communication (e.g., NVLink) bandwidth and inter-node communication (e.g., Infiniband or RoCE) bandwidth since data is communicated on both classes of channels for the NCCL implementation [69]. The exact ratio between the two communication technologies is dependent on factors like the number of nodes and NCCL implementation version (e.g., ring vs. tree). We use real hardware measurement data via to understand what these effective interconnect ratios and bandwidths are in practice. Large-scale training also often exhibits non-constant bandwidth across intra- and inter-node hierarchies. We also consider AllGather and ReduceScatter communication calls, which are required in FSDP and TP.

Computation-Communication Overlap. We maintain separate compute and communication streams and overlap traces with no data dependencies. We also assume GPU kernels are launched whenever data dependencies are resolved. Ideally, we want to maximize compute-communication overlap. However, as demonstrated in Figure 6, there is a segment of exposed communication for the All2All operation, indicating the GPU's compute and memory resources are idle and underutilized.

MAD-Max allows us to both identify combinations of kernels and parallelization strategies that lead to exposed communication and experiment with different parallelization strategies to decrease exposed communication segments. Optimizing for computation-communication overlap is an important objective across multi-node, large-scale ML workloads. Currently, 14~32% of GPU cycles on the training clusters come from exposed communication (Figure 4).

#### V. EXPERIMENTAL METHODOLOGY

This section describes our validation and outlines the design space of this study, including variations of real-world models, hierarchical parallelization strategies, and hardware platforms.

**Performance Model Validation.** Table I lists validation points of various first-order execution metrics across real, measured recommendation and LLM training experiments. For DLRM-A training [40], we validate the performance model over the key dimensions of serialized iteration time,

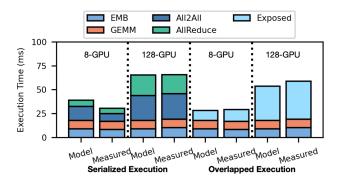


Fig. 7. DLRM-A serialized and overlapped execution validation for 8-, 128-GPU training.

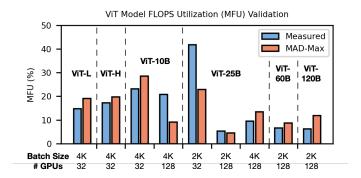


Fig. 8. ViT validation across different model sizes, global batch sizes, and number of GPUs on AWS p4d\_24xlarge instances.

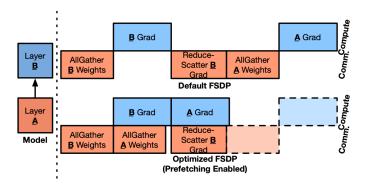


Fig. 9. Optimized FSDP implementation with prefetching that we validate against for production LLaMA training traces.

% communication exposed, and training throughput for 96.89, 91,62, and 99.17% modeling accuracy, respectively. Figure 7 compares the execution time of DLRM-A training in detail across 8- and 128-A100 ZionEX platforms. We validate serialized execution to check layer execution and collectives' volumes, overlapped execution to check at-scale latency-hiding opportunities and systems of different number of nodes to observe networking scaling effects. For DLRM-B training, our model reports 3.05 MQPS whereas the measured throughput is 3.4 MQPS for 89.7% modeling accuracy.

For the largest LLaMA configuration (LLaMA-70B), our performance model estimates training time for all 1.4T tokens

	DLRM-A [40]	DLRM-A Transformer	DLRM-A MoE	DLRM-B [40]	DLRM-B Transformer	DLRM-B MoE	GPT-3 [9]	LLaMA [61]	LLaMA2 [62]	LLM-MoE
# Parameters		793B	795B		332B	333B	175B	65.2B	70B	1.8T
FLOPs per sample/token	638M	2.6B	957M	60M	2.1B	90M	350B	130.4B	140B	550B
Sparse Lookup Bytes per sample/token	22.61 MB		13.19 MB		49.2 KB 32.8 KB		42.8 KB			
Global Batch Size	64K		256K		2K (4M tokens)					
Context Length	N/A	80	N/	A	80	N/A	2	2048	4096	8192

TABLE II
TARGET RECOMMENDATION MODELS, LLMS, AND THEIR VARIANTS BY KEY MODEL-LEVEL CHARACTERISTICS.

	DLRM	LLM	
	Training System [40]	Training System [61]	
Base device	NVIDIA A100 40GB	NVIDIA A100 80GB	
Devices per node	8		
# nodes	16	256	
Peak TF32 throughput	20 PFLOPS	319 PFLOPS	
HBM capacity	5 TB	164 TB	
HBM bandwidth	199 TB/s	3.96 PB/s	
Intra-node interconnect bandwidth (unidirectional)	38.4 TB/s	614.4 TB/s	
Inter-node interconnect fabric	RoCE	Infiniband	
Inter-node interconnect bandwidth (unidirectional)	25.6 Tbps	409.6 Tbps	

TABLE III
BASELINE DISTRIBUTED SYSTEMS USED IN EVALUATION.

to take 19.21 days as opposed to the reported 21 days in [61]. For this use-case, we use the same hardware platform as reported in [61] (i.e., 2048 80GB HBM A100s). We also validate the aggregate GPU Hours to train for 306k steps, resulting in 84.66% modeling accuracy.

Figure 8 presents additional validation points on Vision Transformer (ViT) model training across a range of model configurations, global batch sizes, and number of GPUs. ViT models range from 300M (ViT-L) to 120B (ViT-120B) parameters and global batch size is set at either 2 or 4K for target model accuracies. All experiments are done on AWS p4d\_24xlarge instances and using the baseline FSDP parallelization strategy. We model SM utilization as a function of GPU local batch size and model layer FLOPs requirements. Across all the data points, we get an average of 93.88% and median of 95.74% accuracy for model flops utilization (MFU).

Figure 9 shows a visualization of communication and computation streams for an optimized implementation of FSDP with prefetching enabled. In this optimized variation of FSDP, earlier layer (i.e., Layer A) weight AllGathers are prefetched and overlapped with later layer (i.e., Layer B) gradient computation, leading to overall execution time speedup. We validate this collective-level optimization in MAD-Max against a production implementation and corresponding GPU traces. For a specific LLaMA pre-training run using this optimization, we observe 98% communication overlap against a predicted 93% communication overlap for MAD-Max simulation.

Model Variations. Table II lists the suite of large ML models explored in Section VI. We explore transformer and MoE variants of real-world DLRM-A and DLRM-B. The transformer feature interaction variants have 4 layers and a down-sampled sequence length of 80. MoE variants are configured with 16 experts (2 active) per layer. For the LLM models, we follow specifications in [9], [61], [62]. For LLM-MoE, we explore a hypothetical 1.8T parameter model with

16-(2 active) way MoE for the MLPs in transformer blocks. We use fixed global batch sizes as specified in prior studies [40], [61] to maintain target model accuracy.

**Design Space Exploration.** We use FSDP [75] as the baseline due to its wide adoption and ability to best guarantee training feasibility by minimizing memory footprint. We explore valid hierarchical parallelism strategies at intra- and inter-node levels, considering combinations of DDP, FSDP, and TP. For hardware, unless otherwise stated, we use training systems from prior case studies [40], [61] (Table III). We also explore implications of using H100 and H100 SuperPOD systems by additional simulations replacing our A100-based models with H100 specifications [45], [46] – Table IV.

Validation Efforts. The efforts behind large-scale GPU validation experiments can easily add up. Our DLRM-A and -B validation experiments that were critical for tuning *MAD-Max* (Table I, Figure 7) took ~64K aggregate A100 GPU hours. Running the same experiments on AWS p4d\_24xlarge EC2 instances – which also have 4× lower inter-node interconnect bandwidth compared to systems enumerated in Table III – would amount to even more aggregate GPU hours. Additionally, the LLaMA and ViT validation experiments (Table I, Figure 8), which were run across a range of 32 to 2048 GPUs, would require comparable aggregate GPU hours.

#### VI. EVALUATION RESULTS AND ANALYSIS

When parallelization strategies are tailored to specific deep learning models and tasks at hand, we can achieve 8~124% throughput improvement. Figure 10 overviews pre-training throughput of key large ML models (Table II) normalized to the baseline. We achieve, on average 65.9% pre-training throughput improvement (blue bars) over FSDP by tuning parallelization strategies at the layer-type granularity. The strategy that achieves optimal training throughput is indicated in parenthesis. For example, when considering the base dense layers of DLRM-A, applying Tensor Parallelism within a node of 8 GPUs and Distributed Data Parallelism across nodes of GPUs (i.e., (TP, DDP)) leads to optimal pre-training throughput. In cases like DLRM-A Transformer, where both base dense and transformer layers are present, the optimal way to parallelize each type of layer may differ.

Additionally, the orange dotted bars represent potential throughput improvements from optimizing parallelization strategies without memory constraints of current distributed systems. The optimal parallelization strategy and its expected benefits are influenced by several factors, including the model

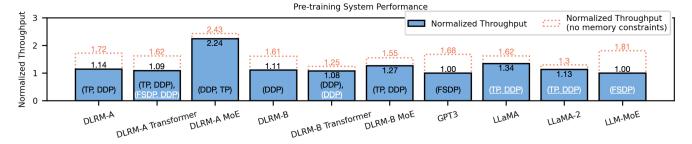


Fig. 10. We can improve pre-training performance over FSDP baseline by applying intra- and inter-node parallelization strategies for base dense and transformer layers separately. Throughput-optimal parallelization strategies are listed in (intra-, inter-) order. Black and white, underlined text refer to recommendation base dense and transformer layers, respectively.

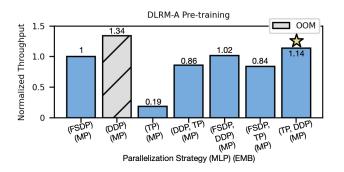


Fig. 11. **DLRM-A Pre-training.** Considering memory capacity constraints, applying TP and DDP for intra- and inter-node parallelism, respectively on base dense layers achieves highest throughput. Gray bar indicates invalid parallelism strategy due to OOM.

architecture, distributed system, and task. We highlight 10 key observations and discuss the underlying insights:

<u>Insight 1: [DLRM]</u> Trillion-parameter embedding tables in <u>DLRMs</u> limit parallelization strategies for the tables to sharding, shifting overall parallelization strategy exploration to focus on the dense components (Figure 11).

Since embedding tables of DLRM-A make up 99.96% of its 793B parameters, the only parallelization strategy viable for DLRM embedding tables on current GPU systems is naive model parallelism sharding. This leaves parallelization strategy exploration on the base dense layers. Figure 11 demonstrates that, over valid parallelization strategies of the base dense layers on the x-axis, training throughput performance of DLRM-A can vary significantly from 0.19 ((TP), (MP)) to 1.14  $\times$  ((TP, DDP), (MP)) over the FSDP baseline. Applying TP scales communication requirements with size of partial sums and activations. If we apply TP at the intra-node level - as opposed to globally - we can take full of advantage of high BW NVLink to communicate the partial sums and activations. In this case, since ((DDP), (MP)) replicates the dense layers' model parameters, gradients, and optimizer states across all devices, causes out-of-memory errors (OOM).

Insight 2: [LLMs] The billion-parameter scale of transformer layers in LLMs makes intra-node replication for compute layers infeasible. In contrast, the small memory footprint of word embeddings (<2GB) allows it to be

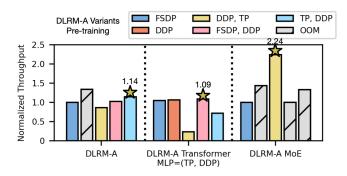


Fig. 12. Between DLRM variants, both optimal parallelization strategy and expected throughput improvement vary.

# replicated across all devices via DDP.

In contrast to DLRMs, for LLMs (e.g., GPT-3), the FSDP offers competitive baseline training throughput (Figure 10). Since the word embeddings of LLMs are relatively small (0.37% of GPT-3), full per-device embedding replication is a viable option via DDP. As in the DLRM cases, we focus our parallelization strategy exploration on the compute-bound layers. However, in the case of GPT-3, any form of layer replication across nodes (e.g., (TP, DDP)) leads to OOM since intra-node sharding is insufficient for meeting memory capacity requirements. Additional device memory capacity can unlock up to 1.68× training throughput improvement.

<u>Insight 3: [Parallelization Strategy Order]</u> Ordering of hierarchical parallelization strategies matter. Replication and sharding strategies must be placed in the correct order to ensure optimal performance. (Figures 10, 11).

The "order" in which we apply hierarchical parallelization strategies matters greatly in terms of both memory capacity footprint and expected throughput. For example, applying ((TP), (DDP)) shards the model component by *number of devices in a node* while applying ((DDP), (TP)) shards the component by *number of nodes*. In Figure 11, where there are 8 GPUs within a node and 16 nodes, the latter strategy leads to a lower per-GPU memory footprint. Additionally, training throughput also varies from using different interconnect channels for communication. For example, ((TP), (DDP)) leads to AllReduce of activations over faster NVLink and weight gradients over slower RoCE/IB. On the other hand,

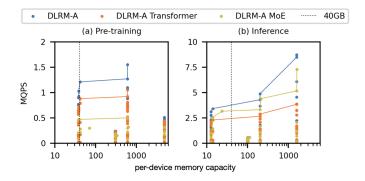


Fig. 13. Pareto curves of parallelization strategies for DLRM variants for (a) pre-training and (b) inference. Each point is a different parallelization strategy.

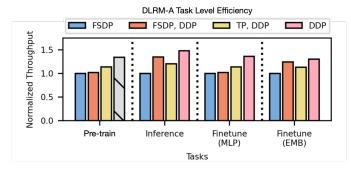


Fig. 14. Task-level diversity (pre-training, inference, and fine-tuning) for the same underlying model and distributed system yields different amounts of speedup over FSDP baselines.

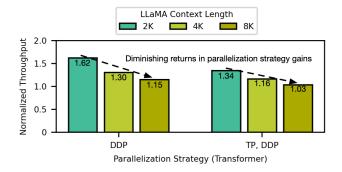


Fig. 15. Given increasing context lengths, solely altering parallelization strategies has diminishing returns for performance benefits over FSDP.

((DDP), (TP)) leads to communicating activations over RoCE/IB and weight gradients over NVLink. For LLMs, long context lengths increase the size of activations to be communicated, so applying inter-node TP leads to significant slowdown (0.18× for GPT-3). On the other hand, utilizing NVLink to communicate large activations leads to  $1.34\times$  speedup for GPT-3.

<u>Insight 4: [DLRM Variants]</u> DLRM Transformer and MoE variants introduce new compute and communication requirements, leading to new parallelization strategy choice and task-level implications. (Figures 12, 13).

Figure 12 shows how the same set of parallelization

strategies interacts with both DLRM-A and its variants. For DLRM-A Transformer, we apply ((TP), (DDP)) on the base dense layers since that is the optimal strategy for DLRM-A and focus parallelization strategy exploration on transformer layers. Across the variants, optimal strategy (yellow star) varies. These differences can be attributed to how transformers introduce more compute and more opportunities for communication-computation overlap while MoE increases blocking, non-overlapping All2All communication. As models continue to evolve, parallelization strategies will as well.

Figure 13 illustrates the parallelization strategy and model architecture options for DLRM-A, categorizing them by required per-device memory and potential throughput for pre-training and inference. The performance-pareto curve is marked with solid lines, indicating that higher memory capacity allows for strategies that achieve greater throughput. For pre-training, the transformer and MoE (Mixture of Experts) variants exhibit lower throughput due to increased computation and communication demands, respectively. During inference, the MoE variant shows greater efficiency compared to the transformer variant as the expensive expert communication is only necessary during the training's backward pass.

<u>Insight 5: [Tasks]</u> Inference, pre-training, and finetuning have different optimal parallelization strategies and scale-out efficiencies due to differences in forward and backward compute graphs (Figure 14).

Figure 14 shows normalized DLRM-A throughput for various parallelization strategies in pre-training, inference, and fine-tuning. For fine-tuning, we also evaluate the two different scenarios of fine-tuning MLP layers and embedding tables.

We see that certain parallelization strategies like DDP may be invalid for pre-training due to their excessive memory footprint requirements from storing per-device replicated model parameters, gradients, and optimizer states. On the contrary, DDP becomes a viable option for inference and fine-tuning since memory footprint requirements are centered around parameters only for inference and parameters with subsets of gradients and optimizer states for fine-tuning. The amount of speedup over FSDP baseline also varies for the different tasks. Fine-tuning exclusively the embedding tables leaves less room for throughput improvement from different MLP sharding strategies. Perhaps counter-intuitively, throughputoptimal parallelization strategy ordering for fine-tuning only embedding tables resembles that for inference. This is because in this scenario we omit the costly MLP weight and input gradient calculations that are found during pre-training.

Insight 6: [Context-Length] Increasing context-lengths limits the improvements from parallelization strategy optimizations, necessitating either changes in model architecture or underlying distributed systems (Figure 15).

Figure 15 shows that input complexity, in terms of context length, plays a key role in training throughput. We investigate the effectiveness of ((DDP)) and ((TP), (DDP)) across LLMs of increasing context lengths. 2K and 4K context length examples refer to LLaMA and LLaMA2 while the 8K context

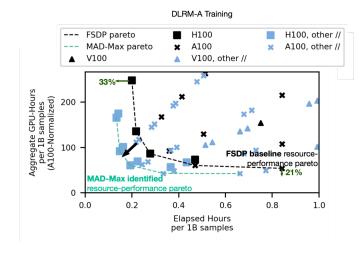


Fig. 16. Across cloud instances of different GPUs and interconnects, parallelization strategy optimizations improve upon the resource-performance pareto frontier of FSDP baseline. Performance is quantified with elapsed time (hr) and compute resource requirements are quantified with aggregate GPU hours (normalized to A100 peak FLOPS).

length data point comes from doubling base LLaMA2's context length while keeping model architecture constant.

We see that throughput gains from tuning parallelization strategy decreases with increasing context length, indicating the limits of optimizing this design space. To further improve throughput performance, changes have to be made to either the underlying distributed system or ML model architecture.

# <u>Insight 7: [Cloud Deployment]</u> Optimizing cloud instance configurations and workload mappings improves both workload performance and operational compute resource requirements.

Figure 16 shows the training time (observed, elapsed hours) and compute resource requirements (aggregate GPU-hours, normalized to A100 peak FLOPS) of training DLRM-A across different GPU-instances from major public cloud providers. To normalize aggregate GPU-hours across different generations of GPUs, we take each experiment's raw aggregate GPUhours and normalize that number by the ratio between the target accelerator's peak FLOPS and A100 peak FLOPS. This normalization is important for reflecting more accurate compute resource requirements since equal amounts of raw, aggregate GPU-hours between two clusters of different compute capabilities should correspond to different levels of resource requirements. For those interested in exploring the tradeoff space for other operational metrics, aggregate GPU-hours can also be potentially converted via other metric-specific ratios. In this example, both performance and resource metrics correspond to processing 1 billion samples – corresponding results for larger workloads (i.e., processing more samples) can be extrapolated using these "per-1B samples" metrics.

The pareto-optimal frontier established from using default FSDP parallelization strategies (black, dotted) can be improved upon by concurrently exploring different instance configurations (number of GPUs, networking capabilities) with par-

	FP-16/32 FLOPS	HBM Capacity, BW	Intra-Node BW (per-device)	Inter-Node BW (per device)
A100 [44]	312, 156 TFLOPS	40GB, 1.6TB/s	600GB/s	200Gbps
H100 [45]	756, 378 TFLOPS	80GB, 2TB/s	900GB/s	400Gbps
H100 SuperPOD [46]	756, 378 TFLOPS	80GB, 2TB/s	900GB/s	1.8TBps
MI250X [3]	383, 96 TFLOPS	128GB, 3.2TB/s	500GB/s	200Gbps
MI300X [5]	1307, 654 TFLOPS	192GB, 5.3TB/s	896GB/s	400GBps
Gaudi2 [24]	400, 200 TFLOPS	96GB, 2.5TB/s	262.5GB/s	300GBps

TABLE IV
SIMULATED COMMODITY HARDWARE SPECIFICATIONS.

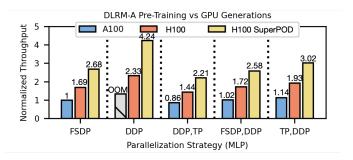


Fig. 17. For DLRM-A pre-training, both overall GPU improvement (H100) and specifically upgrading inter-node interconnect fabric (H100 SuperPOD) lead to observable performance benefits.

allelization strategies (green). As seen in Figure 16's legend, we include three generations of training-class NVIDIA GPUs, ranging from V100s to H100s. For both V100 and A100 instances, both intra- and inter-node interconnect bandwidths vary greatly, with per-device inter-node interconnect bandwidths ranging from <1 to 25GB/s depending on the underlying RoCE or Infiniband specifications. For intra-node interconnect, NVLink-enabled instances provide state-of-theart performance. For this DLRM-A training case study, we see up to 33% training time and 21% compute resource reduction. By extension, operational energy consumption is also reduced due to less compute resources required – as measured by aggregate GPU-hours – for the task at hand. Even if one were to explore this design space via an intelligent, constrained search, having a first-order performance model like MAD-Max for design guidance can enable aggregate GPU-hours savings on the order of 100s per 1 billion samples.

# <u>Insight 8: [GPU-Generations]</u> Across generations of GPUs, improvements in compute, memory, and interconnect not only improve distributed ML performance but also unlock different viable parallelization strategies.

In Figure 17, we compare the A100 against a GPU with H100's specifications via simulation. We also consider the H100 SuperPOD configuration, where the RoCE/IB inter-node interconnect fabric is replaced with NVLink for up to 256 GPUs, leading to  $\sim 4.5 \times$  inter-node interconnect bandwidth compared to H100 DGX systems (see Table IV for full specifications).

Switching from the A100 (blue) to the H100 (orange) results in different levels of performance improvement across various parallelization methods. This variation in speedup is because the enhancements in compute, memory, and networking do not occur at the same rate when upgrading from A100 to H100.

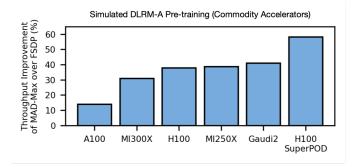


Fig. 18. MAD-Max can simulate other commodity hardware (Table IV) and identify parallelization strategies that improve upon baseline FSDP.

Additionally, each parallelization strategy prioritizes different aspects of the system's resources. Specifically, for DLRM-A training, solely upgrading the inter-node bandwidth (i.e., H100 to H100 SuperPOD) results in 1.82× higher throughput. This is primarily because the inter-node interconnect upgrade directly accelerates the blocking All2All embedding communication collectives.

Insight 9: [Alternative Commodity Hardware] MAD-Max can simulate other commodity hardware platforms with independent compute and communication streams and further identify parallelization strategies with potential performance improvements.

Figure 18 depicts additional simulations for hardware configurations adjusted to best match AMD MI250X, MI300X GPUs and Intel Gaudi2 accelerators (Table IV). Similar to our baseline A100 ZionEX system [40], we evaluate clusters of 128 devices for the DLRM-A pre-training task. For AMD MI GPUs [3], [5], we follow reference scale-out CDNA platform designs [2], [4]. Since Gaudi2 [24] does not have public datasheets, we follow prior benchmarking efforts on Intel Developer Cloud [12]. We show results for throughput improvement from using a MAD-Max identified parallelization strategy over FSDP. Compared to the 40GB-HBM A100, the other hardware platforms' increased HBM capacities (80+GB) allow MAD-Max to identify parallelization strategies that replicate more dense model components for higher pre-training throughput.

Insight 10: [Future Technologies Trends] For large ML workloads, improving individual hardware components leads to limited throughput gain. Unlocking further performance requires jointly improving hardware and systems specifications (Figures 19, 20).

From A100 to H100, compute, memory capacity, memory bandwidth, intra-node interconnect bandwidth, inter-node interconnect bandwidth improve by  $2.42\times, 2\times, 1.29\times, 1.5\times, 2\times$  (9× for SuperPOD), respectively. In Figure 19, we perform a hardware scaling study where compute, memory capacity and bandwidth, intra- and inter-node interconnect bandwidth are all improved by  $10\times$  separately and concurrently. We observe the effects of these improvements on DLRM-A and GPT-3 training and inference.

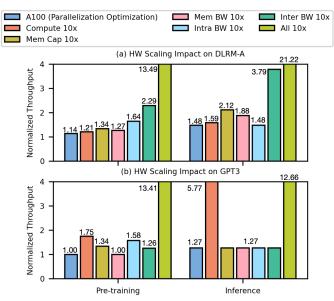


Fig. 19. Individually scaling different hardware capabilities for (a) DLRM-A and (b) GPT-3 workloads leads to sub-linear speedup. Concurrently improving all capabilities leads to super-linear speedup.

For DLRM-A pre-training and inference, independently improving anything but inter-node interconnect by  $10\times$  will only net 1.64 and  $2.12\times$  throughput improvements, respectively. For these use-cases, since blocking All2All embedding communication is performance-critical, targeting inter-node communication bandwidth leads to substantial performance improvement. For GPT-3, since compute-bound layers are critical to overall throughput, improving just compute throughput leads to more workload acceleration compared to DLRMs.

Figure 20 details the sources of the performance changes. Serialized execution breakdown shows execution time allocated to embedding lookups, GEMM, and specific communication collectives, disregarding the effects of overlap. Computation-communication overlap breakdown shows how much communication is hidden behind embedding lookups and GEMM. These breakdowns help us better understand the speedup results from Figure 19 since throughput improvements can come from a variety of sources: accelerating compute-heavy layers (e.g., compute in GPT-3), reducing overall communication time (e.g., All2All in recommendation models), or even unlocking new parallelization strategies with more memory capacity (e.g., DDP for GPT-3).

For all four cases, jointly improving hardware components leads to super-linear performance improvement. This is because distributed ML execution is non-serial so improving the performance of each trace segment can lead to more overlap or unlock new parallelization strategies altogether.

#### VII. RELATED WORK

We discuss related work in two key categories: parallelization strategy and distributed AI performance modeling

**Parallelization Strategy Exploration.** [35], [70] provide compiler annotations for identifying efficient parallelization

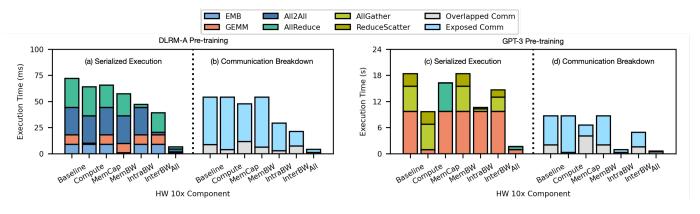


Fig. 20. (a, c) Serialized execution and (b, d) communication breakdown for both DLRM-A and GPT-3 training allows us to better understand where speedups from hardware components come from.

strategies. [37], [56] focus on optimizing communication collectives via fusion and scheduling. [77] focuses on operator-level parallelism. [7], [27] focus on parallelization strategy exploration but are evaluated on older and smaller ML models in Computer Vision and NLP. [73] explores strategies to overlap compute and communication before PyTorch. In this paper, we aim to detach parallelization strategy exploration from existing software implementation details to enable an agile design space exploration of potentially yet to be implemented models. Additionally, we target latest trillion-parameters scale models and expand our design space beyond just collectives.

**Distributed AI Performance Modeling.** [52] provides an analytical model for transformer inference on TPUs. [50] projects computation-communication overlap opportunities for future GPU-centric hardware. [55], [68] provide a simulator for estimating distributed ML performance that is validated against AllReduce collectives. [32] builds upon [55], [68] to introduce a design space exploration tool, yet doesn't focus on optimizing training throughput for specific use cases like DLRM models. These works build upon earlier work in simulating [39], [57] and characterizing [25], [26] distributed systems. [66] emphasizes network optimization. [36] focuses on generating replayable traces to better estimate hardware resource utilization. [60] is an effort to standardize traces across different software implementations for fair comparisons and generating synthetic traces, which can potentially be integrated with our performance model for better integration with current software implementations. We design our performance model to be compatible with different hardware platforms, tasks, and exploration objectives. We also focus on large ML model execution behavior and validate accordingly.

# VIII. CONCLUSION

Training and serving large-scale ML models is a resourceintensive and costly endeavor. We present an agile performance modeling framework to identify efficient solutions for large-scale ML pre-training, fine-tuning, and inference that is also validated against large-scale experiments. Using a suite of real-world large ML models, we identify parallelization strategies for improving performance on existing systems and cloud instances and performance bottlenecks of future hardware systems.

#### ACKNOWLEDGEMENTS

We thank Apostolos Kokolis, Giri Anantharaman, Kalyan Saladi and Srinivas Sridharan for feedback on modeling effective interconnect communication at-scale. We also thank Can Balioglu, Changhan Wang, Kaushik Ram Sadagopan, and Yejin Lee for discussions on performance modeling and optimizations for key deep learning applications.

#### REFERENCES

- B. Acun, M. Murphy, X. Wang, J. Nie, C.-J. Wu, and K. Hazelwood, "Understanding training efficiency of deep learning recommendation models at scale," arXiv preprint arXiv:2011.05497, 2020.
- [2] AMD, "Amd cdna 2 architecture," https://www.amd.com/content/dam/ amd/en/documents/instinct-business-docs/white-papers/amd-cdna2white-paper.pdf, 2021.
- [3] AMD, "Amd instinct mi200 series accelerator," https://www.amd.com/ system/files/documents/amd-instinct-mi200-datasheet.pdf, 2021.
- [4] AMD, "Amd cdna 3 architecture," https://www.amd.com/content/dam/amd/en/documents/instinct-tech-docs/white-papers/amd-cdna-3-white-paper.pdf, 2023.
- [5] AMD, "Amd instinct mi300x accelerator," https://www.amd.com/ content/dam/amd/en/documents/instinct-tech-docs/data-sheets/amdinstinct-mi300x-data-sheet.pdf, 2023.
- [6] R. Anil, S. Gadanho, D. Huang, N. Jacob, Z. Li, D. Lin, T. Phillips, C. Pop, K. Regan, G. I. Shamir, R. Shivanna, and Q. Yan, "On the factory floor: Ml engineering for industrial-scale ads recommendation models," 2022. [Online]. Available: https://arxiv.org/abs/2209.05310
- [7] N. Ardalani, S. Pal, and P. Gupta, "Deepflow: A cross-stack pathfinding framework for distributed ai systems," 2022.
- [8] M. Artetxe, S. Bhosale, N. Goyal, T. Mihaylov, M. Ott, S. Shleifer, X. V. Lin, J. Du, S. Iyer, R. Pasunuru, G. Anantharaman, X. Li, S. Chen, H. Akin, M. Baines, L. Martin, X. Zhou, P. S. Koura, B. O'Horo, J. Wang, L. Zettlemoyer, M. Diab, Z. Kozareva, and V. Stoyanov, "Efficient large scale language modeling with mixtures of experts," in Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics, Dec. 2022, pp. 11699–11732. [Online]. Available: https://aclanthology.org/2022.emnlp-main.804
- [9] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, "Language models are few-shot learners," in Advances in Neural Information Processing Systems, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and

- H. Lin, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 1877–1901. [Online]. Available: https://proceedings.neurips.cc/paper\_files/paper/2020/file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf
- [10] Q. Chen, H. Zhao, W. Li, P. Huang, and W. Ou, "Behavior sequence transformer for e-commerce recommendation in alibaba," in *Proceedings of the 1st International Workshop on Deep Learning Practice for High-Dimensional Sparse Data*, ser. DLP-KDD '19. New York, NY, USA: Association for Computing Machinery, 2019. [Online]. Available: https://doi.org/10.1145/3326937.3341261
- [11] H.-T. Cheng, L. Koc, J. Harmsen, T. Shaked, T. Chandra, H. Aradhye, G. Anderson, G. Corrado, W. Chai, M. Ispir, R. Anil, Z. Haque, L. Hong, V. Jain, X. Liu, and H. Shah, "Wide & deep learning for recommender systems," in *Proceedings of the 1st Workshop on Deep Learning for Recommender Systems*, ser. DLRS 2016. New York, NY, USA: Association for Computing Machinery, 2016, p. 7–10. [Online]. Available: https://doi.org/10.1145/2988450.2988454
- [12] Databricks, "Llm training and inference with intel gaudi 2 ai accelerators," https://www.databricks.com/blog/llm-training-andinference-intel-gaudi2-ai-accelerators, 2024.
- [13] N. Du, Y. Huang, A. M. Dai, S. Tong, D. Lepikhin, Y. Xu, M. Krikun, Y. Zhou, A. W. Yu, O. Firat, B. Zoph, L. Fedus, M. P. Bosma, Z. Zhou, T. Wang, E. Wang, K. Webster, M. Pellat, K. Robinson, K. Meier-Hellstern, T. Duke, L. Dixon, K. Zhang, Q. Le, Y. Wu, Z. Chen, and C. Cui, "GLaM: Efficient scaling of language models with mixture-of-experts," in *Proceedings of the 39th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvari, G. Niu, and S. Sabato, Eds., vol. 162. PMLR, 17–23 Jul 2022, pp. 5547–5569. [Online]. Available: https://proceedings.mlr.press/v162/du22c.html
- [14] A. Firoozshahian, J. Coburn, R. Levenstein, R. Nattoji, A. Kamath, O. Wu, G. Grewal, H. Aepala, B. Jakka, B. Dreyer, A. Hutchin, U. Diril, K. Nair, E. K. Aredestani, M. Schatz, Y. Hao, R. Komuravelli, K. Ho, S. Abu Asal, J. Shajrawi, K. Quinn, N. Sreedhara, P. Kansal, W. Wei, D. Jayaraman, L. Cheng, P. Chopda, E. Wang, A. Bikumandla, A. Karthik Sengottuvel, K. Thottempudi, A. Narasimha, B. Dodds, C. Gao, J. Zhang, M. Al-Sanabani, A. Zehtabioskuie, J. Fix, H. Yu, R. Li, K. Gondkar, J. Montgomery, M. Tsai, S. Dwarakapuram, S. Desai, N. Avidan, P. Ramani, K. Narayanan, A. Mathews, S. Gopal, M. Naumov, V. Rao, K. Noru, H. Reddy, P. Venkatapuram, and A. Bjorlin, "Mtia: First generation silicon targeting meta's recommendation systems," in *Proceedings of the 50th Annual International Symposium on Computer Architecture*, ser. ISCA '23. New York, NY, USA: Association for Computing Machinery, 2023. [Online]. Available: https://doi.org/10.1145/3579371.3589348
- [15] R. Gozalo-Brizuela and E. C. Garrido-Merchan, "Chatgpt is not all you need. a state of the art review of large generative ai models," 2023.
- [16] R. Gozalo-Brizuela and E. C. Garrido-Merchán, "A survey of generative ai applications," 2023.
- [17] U. Gupta, S. Hsia, V. Saraph, X. Wang, B. Reagen, G.-Y. Wei, H.-H. S. Lee, D. Brooks, and C.-J. Wu, "Deeprecsys: A system for optimizing end-to-end at-scale neural recommendation inference," in 2020 ACM/IEEE 47th Annual International Symposium on Computer Architecture (ISCA), 2020, pp. 982–995.
- [18] U. Gupta, C.-J. Wu, X. Wang, M. Naumov, B. Reagen, D. Brooks, B. Cottel, K. Hazelwood, M. Hempstead, B. Jia et al., "The architectural implications of facebook's dnn-based personalized recommendation," in 2020 IEEE International Symposium on High Performance Computer Architecture (HPCA). IEEE, 2020, pp. 488–501.
- [19] K. Hazelwood, S. Bird, D. Brooks, S. Chintala, U. Diril, D. Dzhulgakov, M. Fawzy, B. Jia, Y. Jia, A. Kalro et al., "Applied machine learning at facebook: A datacenter infrastructure perspective," in 2018 IEEE International Symposium on High Performance Computer Architecture (HPCA). IEEE, 2018, pp. 620–629.
- [20] X. He, L. Liao, H. Zhang, L. Nie, X. Hu, and T.-S. Chua, "Neural collaborative filtering," in *Proceedings of the 26th International Conference on World Wide Web*, ser. WWW '17. Republic and Canton of Geneva, Switzerland: International World Wide Web Conferences Steering Committee, 2017, pp. 173–182. [Online]. Available: https://doi.org/10.1145/3038912.3052569
- [21] J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, D. de Las Casas, L. A. Hendricks, J. Welbl, A. Clark, T. Hennigan, E. Noland, K. Millican, G. van den Driessche, B. Damoc, A. Guy, S. Osindero, K. Simonyan, E. Elsen, J. W. Rae, O. Vinyals, and L. Sifre, "Training compute-optimal large language models," 2022.

- [22] S. Hsia, U. Gupta, M. Wilkening, C. Wu, G. Wei, and D. Brooks, "Cross-stack workload characterization of deep recommendation systems," in 2020 IEEE International Symposium on Workload Characterization (IISWC). Los Alamitos, CA, USA: IEEE Computer Society, oct 2020, pp. 157–168. [Online]. Available: https://doi.ieeecomputersociety.org/10.1109/IISWC50251.2020.00024
- [23] S. Hsia, U. Gupta, B. Acun, N. Ardalani, P. Zhong, G.-Y. Wei, D. Brooks, and C.-J. Wu, "Mp-rec: Hardware-software co-design to enable multi-path recommendation," in *Proceedings of the 28th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3*, ser. ASPLOS 2023. New York, NY, USA: Association for Computing Machinery, 2023, p. 449–465. [Online]. Available: https://doi.org/10.1145/3582016.3582068
- [24] Intel, "Intel gaudi2 ai accelerator hl-225b mezzanine card," https:// habana.ai/wp-content/uploads/2023/10/HL-225B\_Datasheet\_10\_23.pdf, 2023
- [25] A. Jain, A. A. Awan, Q. Anthony, H. Subramoni, and D. K. D. Panda, "Performance characterization of dnn training using tensorflow and pytorch on modern clusters," in 2019 IEEE International Conference on Cluster Computing (CLUSTER), 2019, pp. 1–11.
- [26] M. Jeon, S. Venkataraman, A. Phanishayee, u. Qian, W. Xiao, and F. Yang, "Analysis of large-scale multi-tenant gpu clusters for dnn training workloads," in *Proceedings of the 2019 USENIX Conference on Usenix Annual Technical Conference*, ser. USENIX ATC '19. USA: USENIX Association, 2019, p. 947–960.
- [27] Z. Jia, M. Zaharia, and A. Aiken, "Beyond data and model parallelism for deep neural networks," 2018.
- [28] N. P. Jouppi, D. Hyun Yoon, M. Ashcraft, M. Gottscho, T. B. Jablin, G. Kurian, J. Laudon, S. Li, P. Ma, X. Ma, T. Norrie, N. Patil, S. Prasad, C. Young, Z. Zhou, and D. Patterson, "Ten lessons from three generations shaped google's tpuv4i: Industrial product," in 2021 ACM/IEEE 48th Annual International Symposium on Computer Architecture (ISCA), 2021, pp. 1–14.
- [29] N. P. 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. Patterson, "Tpu v4: An optically reconfigurable supercomputer for machine learning with hardware support for embeddings," 2023.
- [30] 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, jun 2020. [Online]. Available: https://doi.org/10.1145/3360307
- [31] N. P. Jouppi, C. Young, N. Patil, D. Patterson, G. Agrawal, R. Bajwa, S. Bates, S. Bhatia, N. Boden, A. Borchers et al., "In-datacenter performance analysis of a tensor processing unit," in Proceedings of the ACM/IEEE 44th Annual International Symposium on Computer Architecture (ISCA). IEEE, 2017, pp. 1–12.
- [32] D. K. Kadiyala, S. Rashidi, T. Heo, A. R. Bambhaniya, T. Krishna, and A. Daglis, "Comet: A comprehensive cluster design methodology for distributed deep learning training," 2022.
- [33] L. Ke, U. Gupta, B. Y. Cho, D. Brooks, V. Chandra, U. Diril, A. Firoozshahian, K. Hazelwood, B. Jia, H.-H. S. Lee et al., "Recnmp: Accelerating personalized recommendation with near-memory processing," in 2020 ACM/IEEE 47th Annual International Symposium on Computer Architecture (ISCA). IEEE, 2020, pp. 790–803.
- [34] Y. Kwon, Y. Lee, and M. Rhu, "Tensordimm: A practical near-memory processing architecture for embeddings and tensor operations in deep learning," in *Proceedings of the 52nd Annual IEEE/ACM International* Symposium on Microarchitecture, 2019, pp. 740–753.
- [35] D. Lepikhin, H. Lee, Y. Xu, D. Chen, O. Firat, Y. Huang, M. Krikun, N. Shazeer, and Z. Chen, "Gshard: Scaling giant models with conditional computation and automatic sharding," 2020.
- [36] M. Liang, W. Fu, L. Feng, Z. Lin, P. Panakanti, S. Zheng, S. Sridharan, and C. Delimitrou, "Mystique: Enabling accurate and scalable generation of production ai benchmarks," in *Proceedings of the 50th Annual International Symposium on Computer Architecture*, ser. ISCA '23. New York, NY, USA: Association for Computing Machinery, 2023. [Online]. Available: https://doi.org/10.1145/3579371.3589072
- [37] K. Mahajan, C.-H. Chu, S. Sridharan, and A. Akella, "Better together: Jointly optimizing ML collective scheduling and execution planning using SYNDICATE," in 20th USENIX Symposium on Networked Systems Design and Implementation (NSDI 23). Boston, MA: USENIX Association, Apr. 2023, pp. 809–824. [Online]. Available: https://www.usenix.org/conference/nsdi23/presentation/mahajan

- [38] J. Manyika, "An overview of bard: an early experiment with generative ai," https://ai.google/static/documents/google-about-bard.pdf, 2023.
- [39] A. Mohammad, U. Darbaz, G. Dozsa, S. Diestelhorst, D. Kim, and N. S. Kim, "dist-gem5: Distributed simulation of computer clusters," in 2017 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS), 2017, pp. 153–162.
- [40] D. Mudigere, Y. Hao, J. Huang, Z. Jia, A. Tulloch, S. Sridharan, X. Liu, M. Ozdal, J. Nie, J. Park, L. Luo, J. A. Yang, L. Gao, D. Ivchenko, A. Basant, Y. Hu, J. Yang, E. K. Ardestani, X. Wang, R. Komuravelli, C.-H. Chu, S. Yilmaz, H. Li, J. Qian, Z. Feng, Y. Ma, J. Yang, E. Wen, H. Li, L. Yang, C. Sun, W. Zhao, D. Melts, K. Dhulipala, K. Kishore, T. Graf, A. Eisenman, K. K. Matam, A. Gangidi, G. J. Chen, M. Krishnan, A. Nayak, K. Nair, B. Muthiah, M. khorashadi, P. Bhattacharya, P. Lapukhov, M. Naumov, A. Mathews, L. Qiao, M. Smelyanskiy, B. Jia, and V. Rao, "Software-hardware co-design for fast and scalable training of deep learning recommendation 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. 993–1011. [Online]. Available: https://doi.org/10.1145/3470496.3533727
- [41] M. Naumov, D. Mudigere, H.-J. M. Shi, J. Huang, N. Sundaraman, J. Park, X. Wang, U. Gupta, C.-J. Wu, A. G. Azzolini *et al.*, "Deep learning recommendation model for personalization and recommendation systems," *arXiv preprint arXiv:1906.00091*, 2019.
- [42] NVIDIA, "Nvidia tesla v100 gpu accelerator datasheet," https://images.nvidia.com/content/technologies/volta/pdf/tesla-voltav100-datasheet-letter-fnl-web.pdf, 2018.
- [43] NVIDIA, "Nvidia dgx-1 deep learning system datasheet," https://www.nvidia.com/content/dam/en-zz/Solutions/Data-Center/dgx-1/dgx-1-rhel-datasheet-nvidia-us-808336-r3-web.pdf, 2019.
- [44] NVIDIA, "Nvidia a100 tensor core gpu datasheet," https://www.nvidia. com/content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-a100-datasheet-us-nvidia-1758950-r4-web.pdf, 2021.
- [45] NVIDIA, "Nvidia h100 tensor core gpu datasheet," https://resources. nvidia.com/en-us-tensor-core/nvidia-tensor-core-gpu-datasheet, 2023.
- [46] NVIDIA, "Nvidia hgx h100 ai supercomputing platform datasheet," https://nvdam.widen.net/s/5kgbjq2v2t/hpc-hgx-h100-datasheet-nvidiaweb, 2023.
- [47] NVIDIA, "Nvidia transformer engine version 0.10.0," https://docs. nvidia.com/deeplearning/transformer-engine/index.html, 2023.
- [48] OpenAI, "Chatgpt," https://chat.openai.com/, 2023.
- [49] OpenAI, "Gpt-4 technical report," 2023.
- [50] S. Pati, S. Aga, M. Islam, N. Jayasena, and M. D. Sinclair, "Computation vs. communication scaling for future transformers on future hardware," 2023
- [51] C. Pei, Y. Zhang, Y. Zhang, F. Sun, X. Lin, H. Sun, J. Wu, P. Jiang, J. Ge, W. Ou, and D. Pei, "Personalized re-ranking for recommendation," in *Proceedings of the 13th ACM Conference on Recommender Systems*, ser. RecSys '19. New York, NY, USA: Association for Computing Machinery, 2019, p. 3–11. [Online]. Available: https://doi.org/10.1145/3298689.3347000
- [52] R. Pope, S. Douglas, A. Chowdhery, J. Devlin, J. Bradbury, A. Levskaya, J. Heek, K. Xiao, S. Agrawal, and J. Dean, "Efficiently scaling transformer inference," 2022.
- [53] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, "Language models are unsupervised multitask learners," 2019.
- [54] S. Rajbhandari, J. Rasley, O. Ruwase, and Y. He, "Zero: Memory optimizations toward training trillion parameter models," 2020.
- [55] S. Rashidi, S. Sridharan, S. Srinivasan, and T. Krishna, "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 (ISPASS), 2020, pp. 81–92.
- [56] 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*, ser. ISCA '22. New York, NY, USA: Association for Computing Machinery, 2022, p. 581–596. [Online]. Available: https://doi.org/10.1145/3470496.3527382
- [57] D. Sanchez and C. Kozyrakis, "Zsim: Fast and accurate microarchitectural simulation of thousand-core systems," in *Proceedings* of the 40th Annual International Symposium on Computer Architecture, ser. ISCA '13. New York, NY, USA: Association for Computing Machinery, 2013, p. 475–486. [Online]. Available: https://doi.org/10.1145/2485922.2485963

- [58] G. Sethi, B. Acun, N. Agarwal, C. Kozyrakis, C. Trippel, and C.-J. Wu, "Recshard: Statistical feature-based memory optimization for industryscale neural recommendation," in 27th ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), 2022.
- [59] M. Shoeybi, M. Patwary, R. Puri, P. LeGresley, J. Casper, and B. Catanzaro, "Megatron-lm: Training multi-billion parameter language models using model parallelism," 2020.
- [60] S. Sridharan, T. Heo, L. Feng, Z. Wang, M. Bergeron, W. Fu, S. Zheng, B. Coutinho, S. Rashidi, C. Man, and T. Krishna, "Chakra: Advancing performance benchmarking and co-design using standardized execution traces," 2023.
- [61] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, A. Rodriguez, A. Joulin, E. Grave, and G. Lample, "Llama: Open and efficient foundation language models," 2023.
- [62] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, D. Bikel, L. Blecher, C. C. Ferrer, M. Chen, G. Cucurull, D. Esiobu, J. Fernandes, J. Fu, W. Fu, B. Fuller, C. Gao, V. Goswami, N. Goyal, A. Hartshorn, S. Hosseini, R. Hou, H. Inan, M. Kardas, V. Kerkez, M. Khabsa, I. Kloumann, A. Korenev, P. S. Koura, M.-A. Lachaux, T. Lavril, J. Lee, D. Liskovich, Y. Lu, Y. Mao, X. Martinet, T. Mihaylov, P. Mishra, I. Molybog, Y. Nie, A. Poulton, J. Reizenstein, R. Rungta, K. Saladi, A. Schelten, R. Silva, E. M. Smith, R. Subramanian, X. E. Tan, B. Tang, R. Taylor, A. Williams, J. X. Kuan, P. Xu, Z. Yan, I. Zarov, Y. Zhang, A. Fan, M. Kambadur, S. Narang, A. Rodriguez, R. Stojnic, S. Edunov, and T. Scialom, "Llama 2: Open foundation and fine-tuned chat models," 2023.
- [63] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin, "Attention is all you need," in Advances in Neural Information Processing Systems, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds., vol. 30. Curran Associates, Inc., 2017. [Online]. Available: https://proceedings.neurips.cc/paper\_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf
- [64] R. Wang, B. Fu, G. Fu, and M. Wang, "Deep & cross network for ad click predictions," 2017.
- [65] R. Wang, R. Shivanna, D. Cheng, S. Jain, D. Lin, L. Hong, and E. Chi, "Dcn v2: Improved deep & cross network and practical lessons for web-scale learning to rank systems," in *Proceedings of the Web Conference 2021*, ser. WWW '21. New York, NY, USA: Association for Computing Machinery, 2021, p. 1785–1797. [Online]. Available: https://doi.org/10.1145/3442381.3450078
- [66] 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," 2022.
- [67] M. Wilkening, U. Gupta, S. Hsia, C. Trippel, C.-J. Wu, D. Brooks, and G.-Y. Wei, "Recssd: Near data processing for solid state drive based recommendation inference," in 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), 2021.
- [68] 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 2023 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS), 2023, pp. 283–294.
- [69] C. Woolley, "Nccl: Accelerated multi-gpu collective communications," https://images.nvidia.com/events/sc15/pdfs/NCCL-Woolley.pdf.
- [70] Y. Xu, H. Lee, D. Chen, B. Hechtman, Y. Huang, R. Joshi, M. Krikun, D. Lepikhin, A. Ly, M. Maggioni, R. Pang, N. Shazeer, S. Wang, T. Wang, Y. Wu, and Z. Chen, "Gspmd: General and scalable parallelization for ml computation graphs," 2021.
- [71] X. Yi, Y.-F. Chen, S. Ramesh, V. Rajashekhar, L. Hong, N. Fiedel, N. Seshadri, L. Heldt, X. Wu, and E. H. Chi, "Factorized deep retrieval and distributed tensorflow serving," in *Proceedings of Machine Learning* and Systems, ser. SysML'18, 2018.
- [72] B. Zhang, L. Luo, X. Liu, J. Li, Z. Chen, W. Zhang, X. Wei, Y. Hao, M. Tsang, W. Wang, Y. Liu, H. Li, Y. Badr, J. Park, J. Yang, D. Mudigere, and E. Wen, "Dhen: A deep and hierarchical ensemble network for large-scale click-through rate prediction," 2022.
- [73] H. Zhang, Z. Zheng, S. Xu, W. Dai, Q. Ho, X. Liang, Z. Hu, J. Wei, P. Xie, and E. P. Xing, "Poseidon: An efficient communication architecture for distributed deep learning on gpu clusters," in *Proceedings of the*

- 2017 USENIX Conference on Usenix Annual Technical Conference, ser. USENIX ATC '17. USA: USENIX Association, 2017, p. 181–193.
- [74] S. Zhang, S. Roller, N. Goyal, M. Artetxe, M. Chen, S. Chen, C. Dewan, M. Diab, X. Li, X. V. Lin, T. Mihaylov, M. Ott, S. Shleifer, K. Shuster, D. Simig, P. S. Koura, A. Sridhar, T. Wang, and L. Zettlemoyer, "Opt: Open pre-trained transformer language models," 2022.
- [75] Y. Zhao, A. Gu, R. Varma, L. Luo, C.-C. Huang, M. Xu, L. Wright, H. Shojanazeri, M. Ott, S. Shleifer, A. Desmaison, C. Balioglu, B. Nguyen, G. Chauhan, Y. Hao, and S. Li, "Pytorch fsdp: Experiences on scaling fully sharded data parallel," 2023.
- [76] Z. Zhao, L. Hong, L. Wei, J. Chen, A. Nath, S. Andrews, A. Kumthekar, M. Sathiamoorthy, X. Yi, and E. Chi, "Recommending what video to watch next: A multitask ranking system," in *Proceedings of the 13th ACM Conference on Recommender Systems*, ser. RecSys '19. New York, NY, USA: ACM, 2019, pp. 43–51. [Online]. Available: http://doi.acm.org/10.1145/3298689.3346997
- [77] L. Zheng, Z. Li, H. Zhang, Y. Zhuang, Z. Chen, Y. Huang, Y. Wang, Y. Xu, D. Zhuo, E. P. Xing, J. E. Gonzalez, and I. Stoica, "Alpa: Automating inter- and Intra-Operator parallelism for distributed deep learning," in 16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22). Carlsbad, CA: USENIX Association, Jul. 2022, pp. 559–578. [Online]. Available: https://www.usenix.org/conference/osdi22/presentation/zheng-lianmin
- [78] G. Zhou, N. Mou, Y. Fan, Q. Pi, W. Bian, C. Zhou, X. Zhu, and K. Gai, "Deep interest evolution network for click-through rate prediction," in Proceedings of the AAAI Conference on Artificial Intelligence, vol. 33, 2019, pp. 5941–5948.
- [79] G. Zhou, X. Zhu, C. Song, Y. Fan, H. Zhu, X. Ma, Y. Yan, J. Jin, H. Li, and K. Gai, "Deep interest network for click-through rate prediction," in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM, 2018, pp. 1059–1068.