

EmbRace: Accelerating Sparse Communication for Distributed Training of Deep Neural Networks

Shengwei Li

National Key Laboratory of Parallel and Distributed Processing
Computer College, National University of Defense Technology
China
lucasleesw9@gmail.com

Zhiquan Lai

National Key Laboratory of Parallel and Distributed Processing
Computer College, National University of Defense Technology
China
zqlai@nudt.edu.cn

Dongsheng Li

National Key Laboratory of Parallel and Distributed Processing
Computer College, National University of Defense Technology
China
dsli@nudt.edu.cn

Yiming Zhang

National Key Laboratory of Parallel and Distributed Processing
Computer College, National University of Defense Technology
China
Xiamen University, China
sdiris@gmail.com

Xiangyu Ye

National Key Laboratory of Parallel and Distributed Processing
Computer College, National University of Defense Technology
China
xyye@nudt.edu.cn

Yabo Duan

National Key Laboratory of Parallel and Distributed Processing
Computer College, National University of Defense Technology
China
yaboduan@nudt.edu.cn

ABSTRACT

Distributed data-parallel training has been widely adopted for deep neural network (DNN) models. Although current deep learning (DL) frameworks scale well for *dense* models like image classification models, we find that these DL frameworks have relatively low scalability for sparse models like natural language processing (NLP) models that have highly sparse embedding tables. Most existing works overlook the sparsity of model parameters thus suffering from significant but unnecessary communication overhead. In this paper, we propose *EmbRace*, an efficient communication framework to accelerate communications of distributed training for sparse models. *EmbRace* introduces *Sparsity-aware Hybrid Communication*, which integrates AlltoAll and model parallelism into data-parallel training, so as to reduce the communication overhead of highly sparse parameters. To effectively overlap sparse communication with both backward and forward computation, *EmbRace* further designs a *2D Communication Scheduling* approach which optimizes the model computation procedure, relaxes the dependency of embeddings, and schedules the sparse communications of each embedding row with a priority queue. We have implemented a prototype of *EmbRace* based on PyTorch and Horovod, and conducted comprehensive evaluations with four representative NLP models. Experimental results show that *EmbRace* achieves up to 2.41× speedup compared to the state-of-the-art distributed training baselines.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ICPP '22, August 29-September 1, 2022, Bordeaux, France

© 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-9733-9/22/08...\$15.00

<https://doi.org/10.1145/3545008.3545011>

CCS CONCEPTS

• **Computing methodologies** → **Parallel computing methodologies; Distributed computing methodologies; Neural networks.**

KEYWORDS

distributed training, deep learning, sparsity of NLP models, communication scheduling

ACM Reference Format:

Shengwei Li, Zhiquan Lai, Dongsheng Li, Yiming Zhang, Xiangyu Ye, and Yabo Duan. 2022. EmbRace: Accelerating Sparse Communication for Distributed Training of Deep Neural Networks. In *51st International Conference on Parallel Processing (ICPP '22)*, August 29-September 1, 2022, Bordeaux, France. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3545008.3545011>

1 INTRODUCTION

Recently, Deep neural networks (DNNs) have been extensively used in different domains like computer vision and natural language processing (NLP). *Data parallelism* [9] on distributed GPUs is the most widely used parallel strategy to reduce the overall training time. Unfortunately, with the fast increasing size of models and number of workers, communication overhead becomes the main performance bottleneck of distributed training [20].

Communication acceleration has been extensively studied in the literature [8, 17, 38]. However, these approaches mainly focus on dense models like DNNs for computer vision. In sparse NLP models, embedding tables are widely-used for feature learning, holding a large portion of parameters. Only a subset of embedding is accessed and exchanged in one training step, bringing dramatic sparsity of embedding parameters in computing and communication, and making it a big challenge to scale the distributed training task efficiently.

When treating all sparse parameters (tensors) as dense, it will involve unnecessary communication overhead and get poor scalability, especially on low-end GPUs. One line of work [12, 21, 36, 43] explores the sparsity-aware communication method, using an individual strategy to aggregate sparse gradients such as AllGather, Model Average and Parameter Server. However, none of these works can reach high efficiency and scalability simultaneously, as the communication overhead of sparse aggregation is still non-negligible.

In addition to communication reduction, recently many works focus on communication scheduling to overlap communication with computation [4, 15, 16, 33]. The key idea is to adjust the communication order of DNN layers to hide more gradient transmission by computation. However, existing scheduling approaches ignore the sparsity of models and the attributes of embedding tables, where model computation depends on the transmission of the entire embedding and the communication of embedding tables becomes the performance bottleneck.

To accelerate communication for distributed training of sparse models, in this paper we propose *EmbRace*, a communication framework that integrates *Sparsity-aware Hybrid Communication* with *2D Communication Scheduling*. Sparsity-aware hybrid Communication applies AlltoAll along with model parallelism and combines AlltoAll and AllReduce together to race the gradients aggregation of sparse models. 2D Communication Scheduling is introduced for effectively overlapping communication and computation by optimizing model computation procedure, employing fine-grained communication scheduling, and relaxing the computation dependency inside of embeddings.

We summarize our contributions as follows:

- We design and implement *EmbRace*, a distributed communication framework for sparse models atop PyTorch [31] and Horovod [36]. We conduct AlltoAll along with column-wise partitioned model parallelism as our communication strategy for sparse embeddings, and analyze its efficiency. To the best of our knowledge, we are the first to bring model-parallel AlltoAll to data-parallel NLP training.
- We propose 2D Communication Scheduling which combines horizontal and vertical scheduling. We dive into embedding matrices, calculate the minimum dependency of the next forward computation and arrange a communication schedule for each row of sparse gradients. This reduces the amount of communication that blocks the FP computation and enables forward computation to overlap more communication. As far as we know, we are the first to realize communication scheduling inside sparse gradients.
- We comprehensively evaluate the performance of *EmbRace* using four representative NLP models on two kinds of GPU clusters with four popular distributed training methods. Experiments show that *EmbRace* accelerates the training process by up to 77.0% on an RTX3090 GPU cluster and up to 141.5% on an RTX2080 GPU cluster, compared to the fastest baseline.

2 BACKGROUND

2.1 Sparsity of NLP Models

A deep learning (DL) model is usually constructed by a sequence of layers and trained by repeating iterations. A typical training step involves three major parts: 1) *forward pass (FP)*, 2) *backward*

Table 1: Model size and embedding size (MB) in popular NLP models. Ratio refers to the proportion of embedding parameters.

Models	Model Size	Embedding Size	Ratio
LM [19]	3186.5	3099.5	97.27%
GNMT-8 [42]	739.1	252.5	34.16%
Transformer [40]	1067.5	263.4	24.67%
BERT-base [10]	417.7	89.4	21.42%

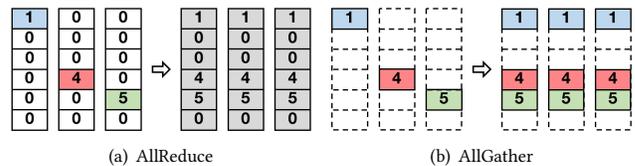


Figure 1: Example sparse data transformation in AllReduce and AllGather with 3 processes.

pass (BP) and 3) *parameter update* with gradients from BP. In data parallelism, gradients are aggregated from all workers after BP to keep model synchronized, where network communication happens.

Most parameters of a DNN model are dense tensors represented by multi-dimensional arrays and stored contiguously in memory. But for sparse tensors full of zeros, only storing the non-zero elements is more efficient in memory footprint. In PyTorch, the coordinate (COO) format is the default sparse data storage format, which describes data as tuples of element indices and their corresponding values.

Embedding matrices are generally adopted in NLP models for feature learning, where each row maps a word to a vector. When training embeddings with batches, a small subset of vocabulary is used, and only the related rows in the embeddings would be updated. Therefore the gradients of embedding matrices could be represented in sparse format, and the communication of embedding gradients could also be sparse. As shown in Table 1, popular NLP models have large portions of embedding parameters. These high proportions of embeddings hint that designing an individual communication scheme for embedding matrices in NLP models is necessary.

2.2 Sparse Communication Architectures

Due to the proven simplicity and high efficiency [24, 37], recent DL frameworks employ collective routines as their communication methods and implement them on top of collective communication libraries such as MPI [13], NCCL [2] and Gloo [1]. In data-parallel training, all workers hold a model replica and synchronize gradients through specific collective primitives before each parameter update operation. Two typical primitives are *AllReduce* and *AllGather*:

- *AllReduce* aggregates data from all processes, reduces the data with an operator such as sum, and distributes results back to each participant. In DNN training, *AllReduce* is widely used for aggregate gradients [36] and implemented effectively [32].

- AllGather gathers the complete data from all workers and distributes the combined data to all workers. AllGather is practical when collecting non-associative tensors such as compressed gradients [27, 39].

AllReduce is the most popular approach when researchers parallelize dense models such as image classification models [17]. However, when distributing DL models which is sparse, AllReduce is no longer a suitable technique. How AllReduce and AllGather communicate sparse data on 3 processes are illustrated in Figure 1, AllReduce has to communicate and sum all data including zeros, while AllGather only sends the non-zero values. Some training methods [18] naively treat the sparse tensors in dense format to utilize the high-performance AllReduce implementation. AllGather is naturally suitable to aggregate the sparse tensors thus popular distributed frameworks [26, 36] employ it for sparse communication. We analyze this further in Section 4.1.2.

2.3 Communication Scheduling

During each BP procedure, the gradients generation depends on the computational graph of DNN. The gradients can start communication immediately after they are calculated rather than wait to complete the entire BP, which is suggested as wait-free backpropagation from Poseidon [44] and supported in most distributed DL frameworks [26, 36].

In DL frameworks, computation and communication are typically carried out with a dependency directed acyclic graph (DAG). In default DAG of DL frameworks, the order of gradients communication is based on a FIFO queue, and FP computations need to wait for the finish of all communications, making communication be overlapped only with the BP computation. Thus comes the *communication scheduling* [16, 33], which switches the FIFO queue into a priority queue to adjust communication orders. The gradient communications that are blocking the beginning of the next FP would be prioritized. So that FP could start earlier and both BP and FP could overlap communication operations.

3 CHALLENGES AND MOTIVATIONS

The sparse and massive embedding tables in NLP models bring serious challenges to efficient communication and scalable distributed training. Previous works pay little attention to this problem, or could not achieve satisfying results. Deeply theoretical and practical analysis of embedding tables motivates us to work on a scalable communication framework for distributed NLP model training.

Efficient communication for embedding tables. Embedding tables occupy an important role in NLP models along with massive parameters, which introduce large communication overhead and seriously restrict the training scalability. Some approaches use hybrid communication architectures combining PS [21] or model average [43] with AllReduce. OmniReduce [12] optimizes AllReduce algorithm to suit the sparse situation. However, based on our analysis in Section 4.1.2 and experience, these solutions are still sub-optimal. Hence, communicating embedding in NLP models with both high efficiency and good scalability, becomes a challenge for us.

Scheduling for sparse communication. Overlapping communication with computation could reduce the communication overhead and accelerate the distributed training. Researchers use

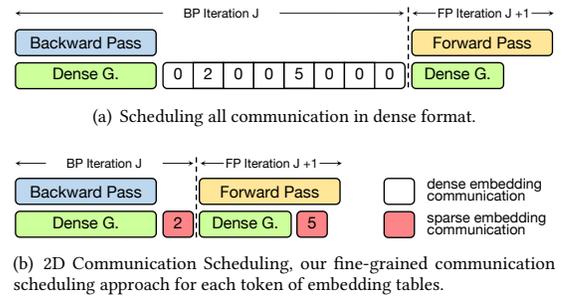


Figure 2: Execution timeline comparisons between communication scheduling approaches, Dense G. refers to the communication operation of dense gradients.

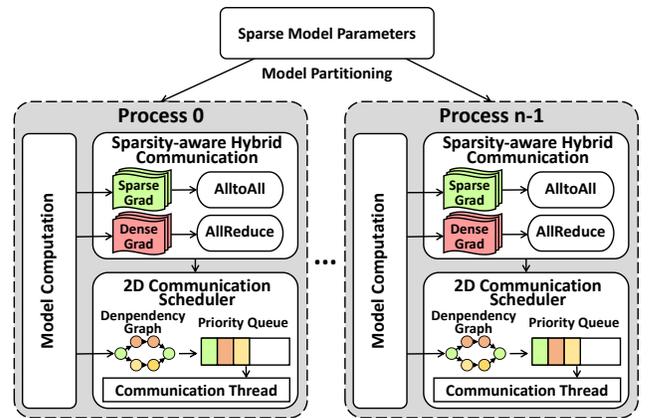


Figure 3: Overview of EmbRace, an efficient communication framework with Sparsity-aware Hybrid Communication and 2D Communication Scheduling.

communication scheduling [4, 15, 16, 33] to achieve this purpose. However, to the best of our knowledge, existing scheduling approaches are only designed for dense models, treating embedding tables as an entire scheduling unit. Embedding tables in dense format will involve unnecessary communication overhead and prevent the starting of forward computation. As shown in Figure 2(a), the communication of zero elements in Embedding would slow the execution of FP. Moreover, embeddings have a unique attribute: parameter updates of embedding tables are element-wise, relating to the tokens which are appeared in a training batch. Based on this property, as illustrated in Figure 2(b), developing an exhaustive scheduling method for NLP tasks to overlap more communication with computation, becomes another challenge for us.

4 EMBRACE DESIGN

Motivated by Section 3, we design EmbRace to address the challenges. Figure 3 shows an overview of EmbRace. To achieve efficient communication for both sparse and dense data in NLP models, we design Sparsity-aware Hybrid Communication, applying hybrid

communication strategies. For sparse tensors, we apply model parallel and AlltoAll to specifically speed up their communications. For dense tensors, AllReduce will take charge as usual. To efficiently schedule the communication of sparse NLP models, we introduce 2D Communication Scheduling. With the help of dependency graph calculation, we could assign priorities to all dense and each block of sparse gradients in model and hold a priority queue to perform all communications operations.

4.1 Sparsity-aware Hybrid Communication

4.1.1 AlltoAll with Model Parallelism. In various NLP tasks, sparse tensors come with embedding tables. As suggested in Section 2.2, there are potential benefits from applying an individual communication scheme for sparse tensors.

Collective operation AlltoAll is useful in recommender system training [30], where processes transmit and receive data slices from every other process to redistribute the data. We bring it with model parallelism to data-parallel NLP training for pursuing the efficiency of sparse communication.

We partition embedding tables for the AlltoAll primitive at first. There are two typical partition solutions, row-wise partitioning and column-wise partitioning. Since each row represents a word in embedding tables, the row-wise partitioning will split different words along with their complete vectors. But the word frequencies are distinct in most datasets, some partitions will be accessed much more frequently, leading to an unbalancing communication cost. In contrast, the column-wise partitioning method will scatter the embedding vectors, keeping a whole corpus vocabulary in each part. The uniform partitions will get the same amount of requests and own balance loads naturally.

The full embedding table is column-wise partitioned into processes before the training start. In each iteration, embedding in each process firstly looks up all training data of this step and produces a different embedding result. Then AlltoAll is called for redistributing the embedding results so that each process gets one embedding result batch and injects it into dense modules. Next, each process computes sparse gradients of embedding tables. Finally, AlltoAll is called again to exchange sparse gradients between processes and each process could update embedding parameters with new gradients.

4.1.2 Efficiency Analyses. Firstly we prove the efficiency of our implement theoretically. Suppose we train an NLP model in a cluster with n nodes, and each node is equipped with w GPU workers. N denotes the total GPU number and $N = w \times n$. There is one embedding table with size M in the model. In each training step, the average density of embedding gradients is α and its sparsity will be $1 - \alpha$, which means the gradient size could be approximated to M when stored as a dense tensor and αM as a sparse tensor. To simplify, we assume that bandwidths B and communication start latency β between any two workers are uniform.

With **AlltoAll**, embedding tables are partitioned among N GPUs. In each training step, AlltoAll is performed twice, one for redistributing embedding lookup results and another for aggregating embedding gradients. A single AlltoAll operation will involve data exchanges with other $N - 1$ training processes, the communication amount of each exchange will be $\frac{\alpha M}{N}$ and the time cost will

be $\frac{\alpha M}{NB} + \beta$. Hence the overall communication overhead will be $2(N - 1)(\frac{\alpha M}{NB} + \beta)$.

AllReduce, which does not support sparse tensor, limits the gradient format to dense. In the most widely used ring-based AllReduce, each GPU divides the gradient into N parts and the communication unit size is $\frac{M}{N}$. Then workers use a ring topology to get the summation of parts by sending and receiving a gradient unit from neighbors for $N - 1$ times. Then, the summed results are gathered through the same ring for another $N - 1$ times. Therefore the overall communication cost is $2(N - 1)(\frac{M}{NB} + \beta)$.

With **Parameter Server (PS)** [25] architecture, we represent the number of servers by S , which holds $S \leq n$. Since we partition the embedding parameters equally across servers, the message size will be $\frac{\alpha M}{S}$ each. During each iteration, servers send parameters to GPUs and receive relative gradients. Both operations transmit message size of $N \times \frac{\alpha M}{S}$, so that the total overhead will be $2N(\frac{\alpha M}{SB} + \beta)$, whose lower bound will be $2N(\frac{\alpha M}{nB} + \beta)$.

In **AllGather**, GPU workers send and receive the gradient in sparse format to each other simultaneously. Since the number of workers is N , the data communication will repeat for $(N - 1)$ times with the gradient size of αM . After summing up the transmitted data size to $(N - 1)\alpha M$, the time will be $(N - 1)(\frac{\alpha M}{B} + \beta)$.

Table 2 summaries the communication costs. Comparing the estimated time results, the AlltoAll, AllReduce, and PS maintain good scalability, where communication will not be strongly impacted by GPU number N . In distributed training environments with $N > 1$ and $N \geq n$, for sparse tensors which hold $\alpha \leq 1$, the AlltoAll method would be faster than AllReduce and PS theoretically. Although in some rare situations with a small number of N and long starting latency β , AlltoAll is slower than AllGather, the transmissions time of AllGather is approximately linear to the GPUs number N with poor scalability.

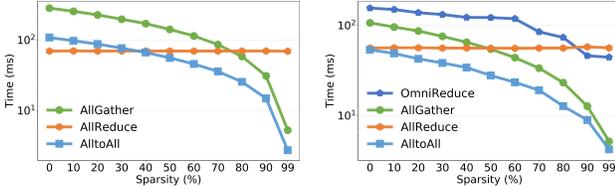
However, in the practical training scenario, different communication algorithms, network topologies and message sizes would influence the bandwidth utilization greatly. Figure 4a shows an example of communication overheads on 8 RTX3090 GPUs, AlltoAll outperforms other methods when the sparsity is greater than 40%. Take models of Table 1 for example, when batch size per worker of LM, GNMT-8, Transformer and BERT-base are 128, 128, 5120 and 32, their corresponding average sparsity are 99.7%, 89.7%, 86.6% and 59.7%. Figure 4b shows another case on 4 RTX3090 GPUs across 4 server nodes, AlltoAll is the best method in all sparsity. OmniReduce could reduce the communication overheads along with the increase of sparsity, but they suffer from insufficient bandwidth usage with excessive divided messages. Therefore, after considering both theoretical and practical results, AlltoAll becomes our sparse communication solution.

4.1.3 Hybrid Communication Architecture. Apart from sparse data, we employ AllReduce architecture to aggregate dense gradients where $\alpha = 1$. The dense parts are treated as an individual dense model so that we could utilize the popular distributed DL framework Horovod.

In summary, we implement a hybrid communication architecture, which combines two collective primitives AllReduce and AlltoAll, and associates data parallelism with model parallelism. We

Table 2: Communication overhead of a sparse tensor according to the communication approaches.

Approaches	Communication Overhead
AlltoAll	$2(N-1)\left(\frac{\alpha M}{NB} + \beta\right)$
AllReduce	$2(N-1)\left(\frac{M}{NB} + \beta\right)$
PS	$2N\left(\frac{\alpha M}{nB} + \beta\right)$
AllGather	$(N-1)\left(\frac{\alpha M}{B} + \beta\right)$



(a) 2 nodes with 4 RTX3090 GPU each (b) 4 nodes with 1 RTX3090 GPU each

Figure 4: Example communication overheads of embedding gradient with different sparsity and communication schemes. Embedding comes from GNMT-8 model and its size is 252.5 MB. OmniReduce [12] implements a sparsity-aware AllReduce algorithm but only supports each node uses 1 GPU.

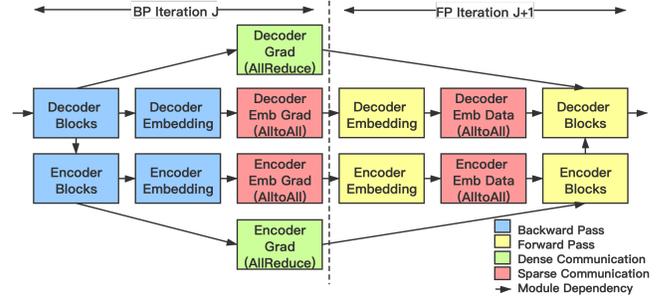
use AlltoAll for sparse communications of embedding tables with column-wise partitioning and prove its efficiency.

4.2 2D Communication Scheduling

To further accelerate distributed training among GPUs, we focus on two main points: reducing the computation stall by overlapping communication with computation; and reducing the communication amount by utilizing the GPUs' idle time. To deeply overlapping, we firstly schedule dense blocks and FP of embedding layers with a priority queue according to the dependency graph. These make our Block-level Horizontal Scheduling. Furthermore, we step into the embedding matrix to reduce the communication quantity. We calculate the minimal gradients dependency, select the necessary gradient rows for the next iteration, prioritize their communication, and delay the remaining communication. These make our Vertical Sparse Scheduling. Combing the horizontal and vertical scheduling, we propose our 2D Communication Scheduling for NLP tasks.

4.2.1 Block-level Horizontal Scheduling. In order to schedule NLP models properly, we explore the computation graph dependency at first. Taking translation models (e.g., Transformer, GNMT-8) as examples, they can be broken down into the following major modules: 1) *Encoder Embedding*; 2) *Encoder Blocks*; 3) *Decoder Embedding*; 4) *Decoder Blocks*. With our Sparsity-aware Hybrid Communication scheme, the dependencies between FP, BP and communication are shown in Figure 5.

In general, the orders of FP and BP are inverse, and the communication would follow a FIFO queue by default. The finish of all gradient communications would prevent starting the FP of the next

**Figure 5: Module dependency graph with Sparsity-aware Hybrid Communication in translation tasks starting from the backward pass. The *Emb Grad* and *Emb Data* refer to the sparse gradients and sparse lookup results of the embedding table.**

iteration, and the FP of Encoder and Decoder Blocks have to wait for their related communicated embedding lookup results. Figure 6(a) depicts the execution timeline of computation and communication with FIFO situation, which is adopted by most distributed frameworks.

To reduce the waiting time, we adopt a priority queue. Unlike popular scheduling methods which use tensor partitioning and involve two inefficiency: extra communication starting overhead due to the increasing number of communication operations; inadequate bandwidth utilization due to the small partitioned message size. We assign the priority according to the following rules. For dense modules of NLP tasks, we observe that the computation loads are more even than image classification models. For example, there are 12 self-attention blocks in Bert-base encoder, each holds a similar number of parameters and takes a comparable calculation time. This similarity hints at applying scheduling with these blocks, where parameters in the same block got the same priority and transmit their gradients together. The dense blocks get priority according to the FP dependency order so that their FP could start as soon as communications finish.

Another characteristic of NLP tasks is that the FP of embedding layers neither depends on each other nor other forward passes. We could perform embedding FP in advance and delay the FP of Encoder Blocks. This running sequence could make room for the AlltoAll communication of embedding data and rescheduled AllReduce communication of dense gradients.

Figure 6(b) illustrates an example execution timeline of Block-level Horizontal Scheduling. Compared to the timeline with default scheduling in Figure 6(a), the communication quantity is the same, but a large amount of communication is overlapped with FP computation. However, limited by bandwidth and module dependency, there is some stalling after BP and some hang between FP.

4.2.2 Vertical Sparse Scheduling. Besides the order adjustment, there are potential reductions in embedding sparse gradients size from the following observations.

Coalescing Gradients. The data batches in NLP are often generated by sentences. On the one hand, there are multiple duplicate words in a single sentence naturally. On the other hand, when a

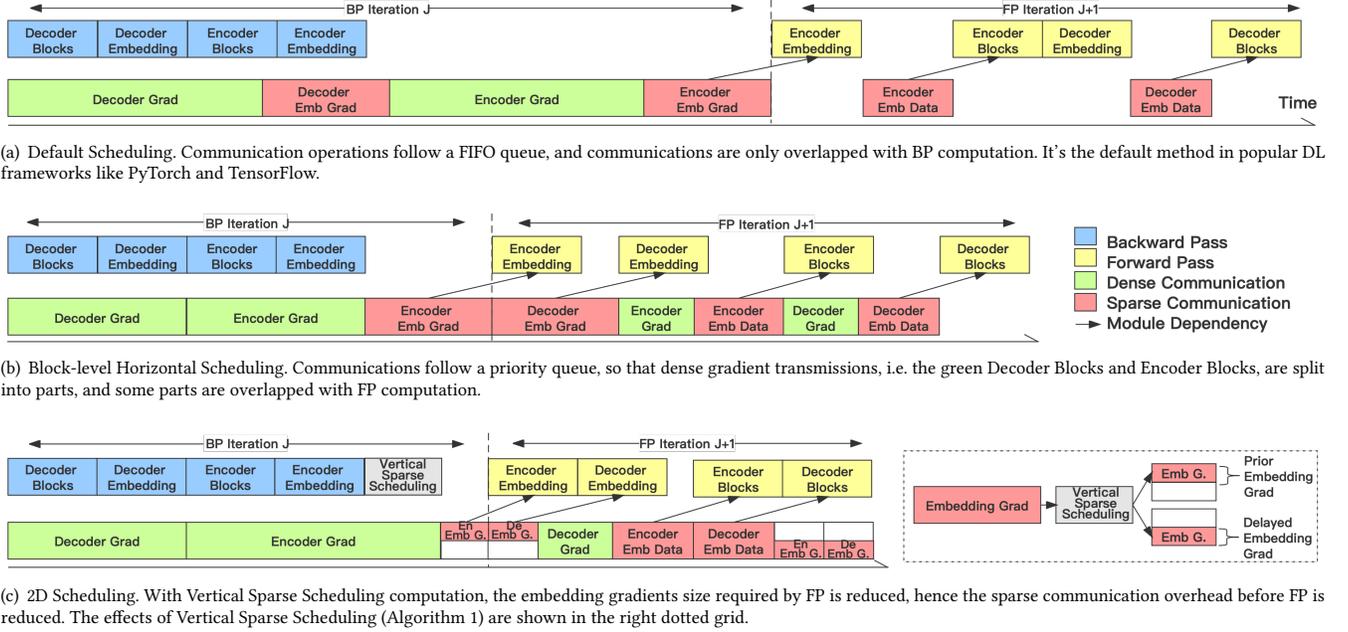


Figure 6: Example execution timelines starting from the backward pass of EmbRace with different communication scheduling schemes. The *Emb Grad (G)* and *Emb Data* refer to the corresponding sparse gradients and sparse lookup results of *Encoder (En)* or *Decoder (De)* Embedding.

tokenizer [23] deals with sentences into uniformly shaped batches, the same value will be padded. With padding and duplicate words, the sparse embedding gradients would have repeated coordinates in the indices. These multi-valued elements could be coalesced into a single value using summation and hence reduce the gradient size. The reduction effect among models is shown in the coalesced size column of Table 3. In LM, GNMT-8, Transformer and BERT-base, the average gradient size is reduced by 20.4%, 53.1%, 52.9% and 84.7%, respectively.

Minimum Dependency of Embedding. Vocabulary is usually much larger than each batch, leading to that only a small subset of embedding is used in each iteration. With the data changing among training steps, different corresponding embedding rows would be updated. This shifting prompts us that the minimum dependency of the subsequent embedding FP is the up-to-date embedding of the following batch data rather than a fully updated

Table 3: Average sparse embedding gradient size (MB) in Vertical Sparse Scheduling. The batch size per worker of LM, GNMT-8, Transformer and BERT-base are 128, 128, 5120 and 32.

Models	Original Grad Size	Coalesced Grad Size	Prioritized Grad Size
LM	8.7	6.9	2.6
GNMT-8	26.0	12.2	5.8
Transformer	35.2	16.6	8.9
BERT-base	36.0	5.5	3.2

embedding. Hence we adopt the data prefetch technology, which always keeps the data of the next iteration in memory. Thanks to the prefetch, we are aware of the data used in the next iteration. When we consider communicating sparse gradients in this iteration, we divide the gradients into two parts, a necessary part and an unhurried part. The necessary part contains the gradients related to the intersection of data between the current and next iteration. In contrast, remaining gradients compose the unhurried part, which could be delayed. Moreover, we assign the highest communication priority to the necessary part and the latest to the unhurried part, making them *prior gradients* and *delayed gradients*. The prioritized gradient size column of Table 3 shows the average size of prior gradients. Compared with the coalesced gradient size, the gradient size is further dropped by 61.8% in LM, 52.5% in GNMT-8, 46.3% in Transformer and 41.9% in BERT-base.

Vertical Scheduling. The cutting down and scheduling will modify the number of rows in sparse embedding gradients. The deciding algorithm is shown in Algorithm 1, line 2-3 represent the computation of coalescing gradients, the line 4-7 refers to a sequence of set operations to distinguish prior gradients. The calculations require a considerable computing resource, and the GPU idle time after BP is a good occasion. The timeline after 2D Communication Scheduling is illustrated in Figure 6(c): the grey box represents the Vertical Sparse Scheduling, whose impact is elucidated at the left dotted grid. The embedding gradient sizes are reduced, and the sparse communication is done in two parts. The communications of prior gradients must be finished before embedding FP and the communications of delayed gradients could be performed later. The embedding FP starts immediately after communicating the

Algorithm 1: Vertical Sparse Scheduling

Input: sparse gradient G , gathered training data for current iteration D_{cur} , gathered training data for next iteration D_{next} , process rank n

Output: prior sparse gradient G_p , delayed sparse gradient G_d

- 1 **After BP:**
- /* Coalesce the duplicate rows */
- 2 $G_{coalesced} \leftarrow COALESCE(G)$
- /* Get the unique training data of process rank n */
- 3 $D_u \leftarrow UNIQUE(D_{cur}[n])$
- 4 $i_{prior} \leftarrow D_u \cap D_{next}$
- 5 $i_{delayed} \leftarrow D_u \setminus i_{prior}$
- /* Select the prior and delayed gradients from coalesced gradients */
- 6 $G_p \leftarrow INDEX_SELECT(G_{coalesced}, i_{prior})$
- 7 $G_d \leftarrow INDEX_SELECT(G_{coalesced}, i_{delayed})$
- 8 return $\{G_p, G_d\}$

slight prior gradients, resulting in better overlapped and efficient data-parallel training.

5 EVALUATIONS

5.1 Implementation

We implement EmbRace on top of DL framework PyTorch 1.8 along with distributed training framework Horovod 0.21.3. EmbRace adopts the open-source communication library NCCL 2.7.8 as the collective routines executor.

EmbRace is integrated with Horovod for usage convenience but takes control of the communication operations. Apart from the same additions of Horovod, EmbRace only requires a few extra lines of code changes in PyTorch to prefetch data and reschedule the FP order. To carry out the Sparsity-aware Hybrid Communication, we first partition the embeddings according to the number of training processes. Then, we replace the AllReduce operations of word embedding gradients with AlltoAll, and we add another AlltoAll to collect embedding FP outputs by registering hooks. As for 2D Communication Scheduling, we hold a priority queue and a communication thread. Communications are performed in the communication thread according to the priority queue. For the Block-level Horizontal Scheduling, before training starts, we assign a priority to each dense block according to the dependency graph and register a hook on each BP of dense blocks. When this hook is fired, the corresponding dense communication operations along with their priorities are dumped into our priority queue. For the Vertical Sparse Scheduling computation, we register another hook on the last BP to finish the calculation, assign corresponding priorities to computation results and add their communications into the priority queue as well.

5.2 Experiments Setup

5.2.1 Hardware Configurations. Two 16-GPU clusters are used in our experimental evaluation. One type of server is equipped with four NVIDIA GeForce RTX3090 GPUs with 24GB GPU memory and six 16G DDR4 RAMs; another has four NVIDIA GeForce RTX2080 GPUs with 8GB GPU memory and three 32G DDR4 RAMs. All

servers have two Intel Xeon 4214R @ 2.40GHz CPUs, run 64-bit Ubuntu 18.04, CUDA 11.1, cuDNN 8.0.5 and are connected by 100 Gbps InfiniBand.

5.2.2 Models and Datasets. We evaluate the performance of EmbRace among four NLP models listed in Table 1. We trained LM model with LM1B [7] dataset, GNMT-8 with WMT-16 En-De [6], Transformer with WMT-14 En-De [5] and BERT-base for question answering task with Squad [34]. The batch size per worker of LM, GNMT-8 and BERT-base are 128, 128 and 32 on RTX3090 GPUs, 128, 32 and 4 on RTX2080 GPUs, respectively. Since the input batch length is flexible in Transformer, we use max tokens per batch of 5120 on RTX3090 GPUs and 500 on RTX2080 GPUs. When measuring the training speed, we use tokens/sec as the metric, where we accumulate the non-padding words in each batch as the number of tokens. All reported throughputs are averaged over the same 500 training steps.

5.2.3 Baselines. Altogether four approaches are compared with EmbRace in our experiments: (i) **BytePS** [18]: a PS based distributed framework that integrates ByteScheduler [33] for partitioning tensors and scheduling communication with a priority queue, but it treats sparse tensors as dense tensors; (ii) **Horovod AllReduce**: the popular distributed communication framework. In the PyTorch implementation of Horovod 0.21.3, the default communication method for sparse tensors is AllReduce; (iii) **Horovod AllGather**: newly introduced with the Horovod 0.22.0 and becomes the default method for sparse data in its PyTorch part, where using AllGather to aggregate sparse tensors and AllReduce for dense tensors; (iv) **Parallax** [21]: a hybrid communication approach that uses a partitioned PS for the sparse communications and AllReduce for the dense communications.

5.3 End-to-End Training Performance

Figure 7 shows the overall training throughputs for LM, GNMT-8, Transformer and BERT-base model on 4, 8 and 16 GPUs. EmbRace achieves a $1.02 \times - 1.77 \times$ speedup on RTX3090 nodes and $1.09 \times - 2.41 \times$ on RTX2080 nodes across four benchmark models. We present the experiment details among different models as follows.

LM: Over the fastest baseline, EmbRace performs $1.18 \times - 1.77 \times$ better on RTX3090 GPUs and $1.99 \times - 2.41 \times$ on RTX2080 GPUs. In our experiments, the LM model has the largest sparse parameter ratio and holds two large embedding tables, each taking over 1.5GB. So that dense communication methods (Horovod AllReduce and BytePs) are too slow. Also limited by the huge embedding tables and GPU memory, RTX3090 GPUs could hold the entire LM model but for RTX2080 GPU we have to put embedding tables on the CPU. For RTX3090 GPUs, the second-fastest method is Horovod AllGather due to its effective NCCL implementation. For RTX2080 GPUs, Horovod AllGather is the second fastest method on 4 GPUs but is surpassed by Parallax on 8 and 16 GPUs, because the scalability of PS is better than AllGather, which is discussed in Section 4.1.2.

GNMT-8/Transformer: EmbRace outperforms the best baseline 10.3%-26.8% in GNMT-8 and 11.5%-17.6% in Transformers with RTX3090 GPUs. With RTX2080 GPUs, EmbRace sees 9.4%-29.6% speedups in GNMT-8 and 10.5%-28.2% speedups in Transformers. There are fewer sparse parameters than LM in translation

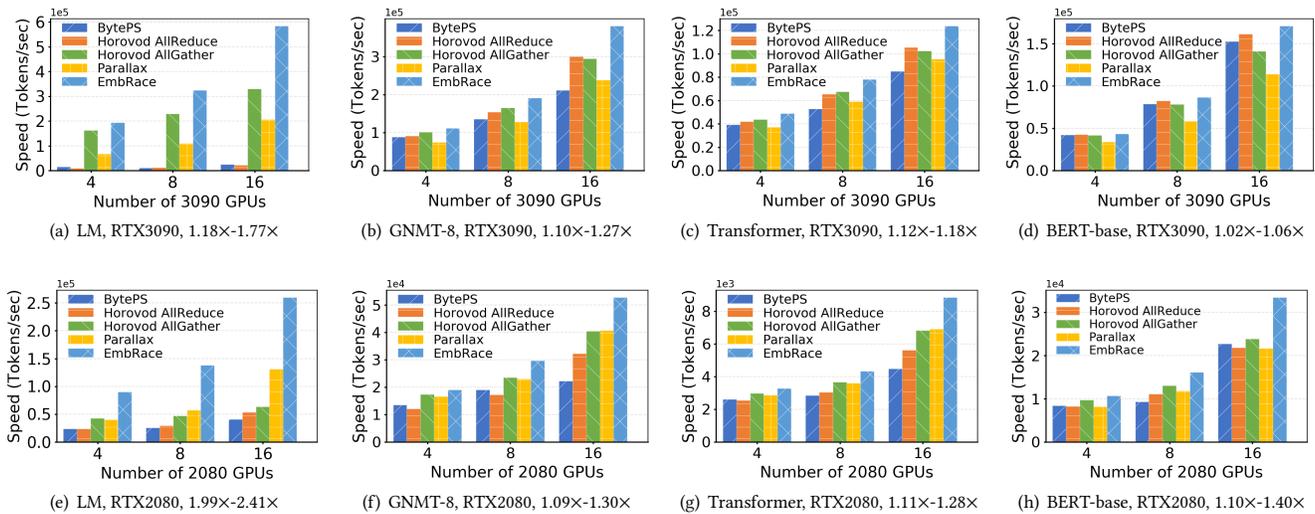


Figure 7: The training performance of NLP models using different numbers of GPUs. The numbers are speedups of EmbRace over the best baseline.

models GNMT-8 and Transformer. When training GNMT-8 and Transformer with RTX3090 GPUs, Horovod AllGather is the fastest method on 4 and 8 GPUs, but performs worse than Horovod AllReduce on 16 GPUs due to its scalability. BytePs performs worse than Horovod AllReduce because BytePs uses share memory to speed up communication. In our hardware environment, the speed of RAMs is slow and would damage the performance of BytePs. The Parallax is less effective than Horovod AllReduce due to the frequent memory copy between GPU and CPU. However on RTX2080 GPUs, dense methods get poor performance due to the smaller batch size and lower intra-node bandwidth. Horovod AllGather is the fastest baseline on 4 and 8 GPUs while exceeded by Parallax on 16 GPUs.

BERT-base: When compared to the best baseline, EmbRace achieves 1.02×-1.06× throughput on RTX3090 GPUs and 1.10×-1.40× speedups on RTX2080 GPUs. BERT-base holds the least number of parameters and the minimum ratio of computation to communication across benchmark models. In RTX3090 clusters, the BP process is long enough for Horovod AllReduce and BytePs to transmit a dense formatted embedding table. Hence Horovod AllReduce and BytePs become faster and 2D Communication Scheduling provides less performance improvement. As for RTX2080 clusters, the communication becomes an obvious bottleneck with the decreasing batch size thus EmbRace gains more accelerations.

5.4 Communication Efficiency

As EmbRace focuses on accelerating sparse communication in model training, the end-to-end training speed results might not be enough. We bring out another index, *Computation Stall*, to further evaluate communication efficiency. In this paper, the Computation Stall is defined as the computation stall time caused by communication during the training procedure. For EmbRace, the Computation Stall consists of the Vertical Sparse Scheduling computation and communications that are not overlapped by computation. While for

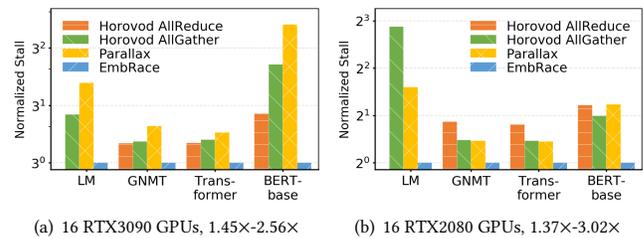


Figure 8: Computation Stall comparing of LM, GNMT-8, Transformer and BERT. Stall values are normalized by EmbRace.

other approaches, Computation Stall is only the non-overlapping communication overheads.

Computation Stall is mainly affected by three conditions: the network bandwidth, communication efficiency, and overlap ratio. Predictably, higher bandwidth and communication efficiency lead to a lower Computation Stall. With more unsatisfactory GPU performance or larger batch size, the computation time becomes longer and thus overlap ratio becomes higher, reducing the Computation Stall as well. In EmbRace, Sparsity-aware Hybrid Communication concentrates on boosting each sparse communication operation and 2D Communication Scheduling aims to overlap more communication with computation. We expect that these two techniques work together and reduce the Computation Stall greatly.

Figure 8 shows the Computation Stall comparison, whose values are normalized by EmbRace, among four benchmark models on 16 GPUs. In LM, Horovod AllReduce gets such a large Computation Stall that it is hard to illustrate, so we omit it. EmbRace sees 1.45×-2.56× speedups on 16 RTX3090 GPUs and 1.37×-3.02× on 16 RTX2080 GPUs. EmbRace performs surprisingly in LM and

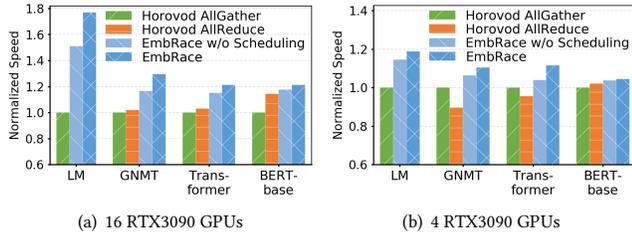


Figure 9: Ablation experiments on the optimizations of EmbRace with 16 and 4 RTX3090 GPUs. Training speeds are normalized by Horovod AllGather.

BERT-base, providing at least 49.7% reduction of Computation Stall. For GNMT-8 and Transformer, although the large amount of dense communication limits further accelerations, EmbRace does a good job in sparse communication and scheduling, cutting down a minimum of 26.4% Computation Stall. The reduced Computation stalls lead to the end-to-end speedup of EmbRace. For example, when training BERT-base on RTX3090 GPUs, the batch size is larger, and the BP time can cover most communications. So computation stalls account for a low proportion of the overall time. EmbRace obviously reduces computation stalls, but the speedups in end-to-end performance are minor. While on RTX2080 GPUs, the computation stall becomes a major factor affecting training performance. The stall reductions of EmbRace bring more speedups to end-to-end performance.

5.5 Ablation Study

To analyze the effectiveness of our optimizations, we carry out the ablation study with all four benchmark models on RTX3090 GPUs. When compared with Horovod AllGather and Horovod AllReduce, results of EmbRace without Scheduling could show the performance gains of Sparsity-aware Hybrid Communication. To demonstrate the effect of 2D Communication Scheduling, we could compare the training speeds of EmbRace with and without Scheduling.

Figure 9 presents the experiment results, where training speed values are normalized by Horovod AllGather. Sparsity-aware Hybrid Communication provides 2.9%-51.0% speedups and 2D Communication Scheduling brings another 3.0%-26.0% speedups on 16 GPUs. On 4 GPUs, Sparsity-aware Hybrid Communication offers 1.5%-14.6% speedups and 2D Communication Scheduling sees another 0.7%-7.5% speedups. Both techniques of EmbRace could bring remarkable improvements. With the increasing number of GPUs, communication accelerations become more obvious.

5.6 Scalability

To demonstrate the scalability of EmbRace, we compare it with Horovod AllReduce which owns the best scalability in GNMT-8, Transformer and BERT-base. For the LM model, we choose Parallax as the competitor since dense methods are extremely slow, and PS architecture maintains good scalability in our analysis.

Figure 10 shows the training throughputs compared with corresponding ideal linear scaling results. We use the related throughputs of models on 4 RTX3090 GPUs as the benchmark, then calculate

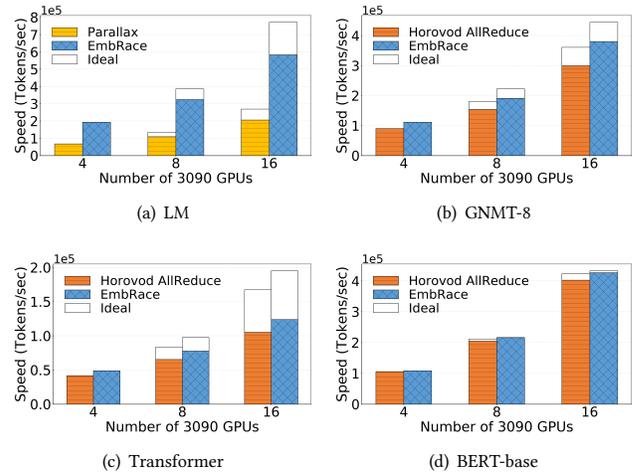


Figure 10: Scaling performance on RTX3090 GPUs, compared to the approach with the second-best scalability.

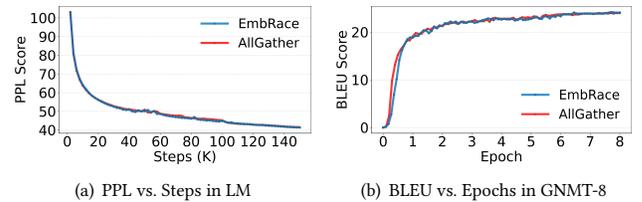


Figure 11: Convergence comparing between EmbRace and Horovod AllGather on 8 RTX3090 GPUs.

the ideal linear results of 8 and 16 GPUs based on the benchmark. When scaled from 4 GPUs to 16 GPUs across GNMT-8, Transformer and BERT-base, EmbRace offers 3.42 \times , 2.53 \times and 3.94 \times speedups while Horovod AllReduce gets 3.32 \times , 2.51 \times and 3.81 \times , respectively. When scaling LM to 16 GPUs, EmbRace gets 3.14 \times throughput, where Parallax gets 3.06 \times . Due to the limitation of the number of devices, we did not test on more server nodes. With the help of better scalability, we expect that EmbRace will have more significant advantages on more GPUs.

5.7 Convergence

Since EmbRace is a synchronous data-parallel training approach, it should correctly converge models like other synchronous approaches. However with Vertical Sparse Scheduling, each sparse gradient is divided into two parts and requires parameter updating twice with sparse optimizers. Because the common sparse optimizer such as Adagrad [11] and SGD [35] is fully element-wise, no matter whether updating embedding matrices with multiple gradient parts or a whole gradient would lead to the same result. But we used Adam [22] optimizer in all experiments. Most parts of Adam are element-wise except the state parameter *step*, which would be accumulated after every optimizing call, leading to a minor difference between a single or two parameter updates. Therefore, we modify

the Adam optimizer in PyTorch, updating the *step* state only at applying the delayed sparse gradients to embedding parameters. This modification ensures synchronous training and the rate of convergence.

Figure 11 illustrates the convergence curves of EmbRace and Horovod AllGather in LM and GNMT-8 with 8 RTX3090 GPUs. We trace the perplexity (PPL) scores for LM with training steps and BLEU scores for GNMT-8 along with epoch numbers. As shown in Figure 11(a) and Figure 11(b), both methods converge the model into PPL 41.5 and BLEU 24.0 in similar numbers of training iterations or epochs.

6 RELATED WORK

Communication acceleration. Upon popular synchronous data-parallel frameworks such as Horovod [36] and PyTorch Distributed [26], there are different directions for accelerating communications: (1) optimizing multiple GPU intra-machine communication [41]; (2) applying topology-aware hierarchical collective communication [8, 29]; (3) reducing messages size with gradient compression [3, 27]; (4) speeding up individual messages with specific networks such as RDMA [28]. These works concentrate on each communication operation and could be orthogonal and complementary to EmbRace.

Sparse communication in DNN training. To improve sparse communication performance in DNN training, Horovod and PyTorch Distributed adopt AllGather for sparse communications; Parallax [21] uses architecture combining PS with AllReduce; HMA [43] integrates model average with AllReduce; OmniReduce [12] implements a sparsity-aware AllReduce algorithm; and S2 reducer [14] designs a compress algorithm that adapts sparse tensors to AllReduce primitives. However, they do not utilize the property of embedding matrix, which could reduce the communication overhead further.

Communication scheduling. On top of the wait-free backward propagation [44] which is supported by most DL frameworks, communication scheduling could further overlap communication with FP computation. TicTac [15] attempts to schedule communication in PS architecture with a priority queue; P3 [16] proposes a similar idea and implements layer partitioning for scheduling in a layer-wise granularity; ByteScheduler [33] expands communication scheduling to AllReduce architecture and adopts parameter partitioning for a finer granularity; PACE [4] changes the priority queue into a preemptive queue in AllReduce and uses tensor fusion for better bandwidth usage. However, these methods only work for dense data. In EmbRace, we benefit from these approaches but we take sparsity and characteristics of NLP models into consideration.

7 CONCLUSION

In this paper, we present EmbRace, an efficient distributed sparse communication framework for NLP model training. EmbRace introduces Sparsity-aware Hybrid Communication that uses AlltoAll primitive with model parallelism to race the communication of embedding tables. We also design 2D Communication Scheduling which embraces the Horizontal and Vertical Scheduling, letting EmbRace optimizes model computation procedure, utilizes GPU idle time, and achieves a thorough overlapped communication schedule. Experiments show EmbRace achieves up to 66.7% Computation

Stall reduction, up to 2.41× training speedup, and better scalability, when compared to the best baseline. Although training giant NLP models with thousands of GPUs becomes a trend nowadays, training models swiftly with limited resources such as RTX2080 GPUs still matter. And EmbRace could benefit sparse communications in giant NLP models training as well.

ACKNOWLEDGMENTS

This work is sponsored by National Key R&D Program of China (2021YFB0301200). Zhiqian Lai is the corresponding author of this paper.

REFERENCES

- [1] 2019. Gloo: a collective communications library. <https://github.com/facebookincubator/gloo>
- [2] 2019. NVIDIA Collective Communications Library (NCCL). <https://developer.nvidia.com/nccl>
- [3] Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. 2017. QSGD: Communication-efficient SGD via gradient quantization and encoding. *Advances in Neural Information Processing Systems* 30 (2017), 1709–1720.
- [4] Yixin Bao, Yanghua Peng, Yangrui Chen, and Chuan Wu. 2020. Preemptive all-reduce scheduling for expediting distributed dnn training. In *IEEE INFOCOM 2020-IEEE Conference on Computer Communications*. IEEE, 626–635.
- [5] Ondřej Bojar, Christian Buck, Christian Federmann, Barry Haddow, Philipp Koehn, Johannes Leveling, Christof Monz, Pavel Pecina, Matt Post, Herve Saint-Amand, et al. 2014. Findings of the 2014 workshop on statistical machine translation. In *Proceedings of the ninth workshop on statistical machine translation*. 12–58.
- [6] Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Varvara Logacheva, Christof Monz, et al. 2016. Findings of the 2016 conference on machine translation. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*. 131–198.
- [7] Ciprian Chelba, Tomas Mikolov, Mike Schuster, Qi Ge, Thorsten Brants, Philipp Koehn, and Tony Robinson. 2013. One billion word benchmark for measuring progress in statistical language modeling. *arXiv preprint arXiv:1312.3005* (2013).
- [8] Minsik Cho, Ulrich Finkler, and David Kung. 2019. BlueConnect: Novel hierarchical all-reduce on multi-tired network for deep learning. In *Proceedings of the 2nd SysML Conference*.
- [9] Jeffrey Dean, Greg S Corrado, Rajat Monga, Kai Chen, Matthieu Devin, Quoc V Le, Mark Z Mao, Marc'Aurelio Ranzato, Andrew Senior, Paul Tucker, et al. 2012. Large scale distributed deep networks. (2012).
- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [11] John Duchi, Elad Hazan, and Yoram Singer. 2011. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research* 12, 7 (2011).
- [12] Jiawei Fei, Chen-Yu Ho, Atal Narayan Sahu, Marco Canini, and Amedeo Sapio. 2020. *Efficient Sparse Collective Communication and its application to Accelerate Distributed Deep Learning*. Technical Report.
- [13] Edgar Gabriel, Graham E Fagg, George Bosilca, Thara Angskun, Jack J Dongarra, Jeffrey M Squyres, Vishal Sahay, Prabhajan Kambadur, Brian Barrett, Andrew Lumsdaine, et al. 2004. Open MPI: Goals, concept, and design of a next generation MPI implementation. In *European Parallel Virtual Machine/Message Passing Interface Users' Group Meeting*. Springer, 97–104.
- [14] Keshi Ge, Yongquan Fu, Yiming Zhang, Zhiqian Lai, Xiaoge Deng, and Dongsheng Li. 2022. S2 reducer: High-performance sparse communication to accelerate distributed deep learning. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 5233–5237.
- [15] Sayed Hadi Hashemi, Sangeetha Abdu Jyothi, and Roy H Campbell. 2018. Tictac: Accelerating distributed deep learning with communication scheduling. *arXiv preprint arXiv:1803.03288* (2018).
- [16] Anand Jayarajan, Jinliang Wei, Garth Gibson, Alexandra Fedorova, and Gennady Pekhimenko. 2019. Priority-based parameter propagation for distributed DNN training. *arXiv preprint arXiv:1905.03960* (2019).
- [17] Xianyan Jia, Shutao Song, Wei He, Yangzihao Wang, Haidong Rong, Feihu Zhou, Liqiang Xie, Zhenyu Guo, Yuanzhou Yang, Liwei Yu, et al. 2018. Highly scalable deep learning training system with mixed-precision: Training imagenet in four minutes. *arXiv preprint arXiv:1807.11205* (2018).
- [18] Yimin Jiang, Yibo Zhu, Chang Lan, Bairen Yi, Yong Cui, and Chuanxiong Guo. 2020. A Unified Architecture for Accelerating Distributed DNN Training in

- Heterogeneous GPU/CPU Clusters. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*. 463–479.
- [19] Rafal Jozefowicz, Oriol Vinyals, Mike Schuster, Noam Shazeer, and Yonghui Wu. 2016. Exploring the limits of language modeling. *arXiv preprint arXiv:1602.02410* (2016).
- [20] Janis Keuper and Franz-Josef Preundt. 2016. Distributed training of deep neural networks: Theoretical and practical limits of parallel scalability. In *2016 2nd Workshop on Machine Learning in HPC Environments (MLHPC)*. IEEE, 19–26.
- [21] Soojeong Kim, Gyeong-In Yu, Hoin Park, Sungwoo Cho, Eunji Jeong, Hyeonmin Ha, Sanha Lee, Joo Seong Jeong, and Byung-Gon Chun. 2019. Parallax: Sparsity-aware data parallel training of deep neural networks. In *Proceedings of the Fourteenth EuroSys Conference 2019*. 1–15.
- [22] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [23] Taku Kudo and John Richardson. 2018. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. *arXiv preprint arXiv:1808.06226* (2018).
- [24] Dongsheng Li, Zhiqian Lai, Keshi Ge, Yiming Zhang, Zhaoning Zhang, Qinglin Wang, and Huaimin Wang. 2019. HPDL: towards a general framework for high-performance distributed deep learning. In *2019 IEEE 39th International Conference on Distributed Computing Systems (ICDCS)*. IEEE, 1742–1753.
- [25] Mu Li, David G Andersen, Jun Woo Park, Alexander J Smola, Amr Ahmed, Vanja Josifovski, James Long, Eugene J Shekita, and Bor-Yiing Su. 2014. Scaling distributed machine learning with the parameter server. In *11th USENIX Symposium on Operating Systems Design and Implementation (OSDI 14)*. 583–598.
- [26] Shen Li, Yanli Zhao, Rohan Varma, Omkar Salpekar, Pieter Noordhuis, Teng Li, Adam Paszke, Jeff Smith, Brian Vaughan, Pritam Damania, et al. 2020. Pytorch distributed: Experiences on accelerating data parallel training. *arXiv preprint arXiv:2006.15704* (2020).
- [27] Yujun Lin, Song Han, Huizi Mao, Yu Wang, and William J Dally. 2017. Deep gradient compression: Reducing the communication bandwidth for distributed training. *arXiv preprint arXiv:1712.01887* (2017).
- [28] Jiuxing Liu, Jiesheng Wu, and Dhabaleswar K Panda. 2004. High performance RDMA-based MPI implementation over InfiniBand. *International Journal of Parallel Programming* 32, 3 (2004), 167–198.
- [29] Hiroaki Mikami, Hisahiro Sugauma, Yoshiki Tanaka, Yuichi Kageyama, et al. 2018. Massively distributed SGD: ImageNet/ResNet-50 training in a flash. *arXiv preprint arXiv:1811.05233* (2018).
- [30] Dheevatsa Mudigere, Yuchen Hao, Jianyu Huang, Andrew Tulloch, Srinivas Sridharan, Xing Liu, Mustafa Ozdal, Jade Nie, Jongsoo Park, Liang Luo, et al. 2021. High-performance, distributed training of large-scale deep learning recommendation models. *arXiv preprint arXiv:2104.05158* (2021).
- [31] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* 32 (2019), 8026–8037.
- [32] Pitch Patarasuk and Xin Yuan. 2009. Bandwidth optimal all-reduce algorithms for clusters of workstations. *J. Parallel and Distrib. Comput.* 69, 2 (2009), 117–124.
- [33] Yanghua Peng, Yibo Zhu, Yangrui Chen, Yixin Bao, Bairen Yi, Chang Lan, Chuan Wu, and Chuanxiong Guo. 2019. A generic communication scheduler for distributed dnn training acceleration. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles*. 16–29.
- [34] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250* (2016).
- [35] Sebastian Ruder. 2016. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747* (2016).
- [36] Alexander Sergeev and Mike Del Balso. 2018. Horovod: fast and easy distributed deep learning in TensorFlow. *arXiv preprint arXiv:1802.05799* (2018).
- [37] Shaohuai Shi, Qiang Wang, and Xiaowen Chu. 2018. Performance modeling and evaluation of distributed deep learning frameworks on gpus. In *2018 IEEE 16th Intl Conf on Dependable, Autonomic and Secure Computing, 16th Intl Conf on Pervasive Intelligence and Computing, 4th Intl Conf on Big Data Intelligence and Computing and Cyber Science and Technology Congress (DASC/PiCom/DataCom/CyberSciTech)*. IEEE, 949–957.
- [38] Peng Sun, Wansen Feng, Ruobing Han, Shengen Yan, and Yonggang Wen. 2019. Optimizing network performance for distributed dnn training on gpu clusters: Imagenet/alexnet training in 1.5 minutes. *arXiv preprint arXiv:1902.06855* (2019).
- [39] Jesper Larsson Träff, Andreas Ripke, Christian Siebert, Pavan Balaji, Rajeev Thakur, and William Gropp. 2008. A simple, pipelined algorithm for large, irregular all-gather problems. In *European Parallel Virtual Machine/Message Passing Interface Users' Group Meeting*. Springer, 84–93.
- [40] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*. 5998–6008.
- [41] Guanhua Wang, Shivaram Venkataraman, Amar Phanishayee, Jorgen Thelin, Nikhil Devanur, and Ion Stoica. 2019. Blink: Fast and generic collectives for distributed ml. *arXiv preprint arXiv:1910.04940* (2019).
- [42] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. 2016. Google's neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144* (2016).
- [43] Yuetong Yang, Zhiqian Lai, Lei Cai, and Dongsheng Li. 2020. Model Average-based Distributed Training for Sparse Deep Neural Networks. In *IEEE INFOCOM 2020-IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS)*. IEEE, 1346–1347.
- [44] Hao Zhang, Zeyu Zheng, Shizhen Xu, Wei Dai, Qirong Ho, Xiaodan Liang, Zhiting Hu, Jinliang Wei, Pengtao Xie, and Eric P Xing. 2017. Poseidon: An efficient communication architecture for distributed deep learning on GPU clusters. In *2017 USENIX Annual Technical Conference (USENIX ATC 17)*. 181–193.