# Optimizing Deep Learning Recommender Systems' Training On CPU Cluster Architectures

Dhiraj Kalamkar*, Evangelos Georganas†, Sudarshan Srinivasan*,
Jianping Chen‡, Mikhail Shiryaev§, and Alexander Heinecke†
*Intel Technology India Private Ltd., India
†Intel Corporation, USA
‡Intel China Co. Ltd., China
§Intel Russia, Russia

*Abstract*—During the last two years, the goal of many researchers has been to squeeze the last bit of performance out of HPC system for AI tasks. Often this discussion is held in the context of how fast ResNet50 can be trained. Unfortunately, ResNet50 is no longer a representative workload in 2020. Thus, we focus on Recommender Systems which account for most of the AI cycles in cloud computing centers. More specifically, we focus on Facebook's DLRM benchmark. By enabling it to run on latest CPU hardware and software tailored for HPC, we are able to achieve more than two-orders of magnitude improvement in performance (110x) on a single socket compared to the reference CPU implementation, and high scaling efficiency up to 64 sockets, while fitting ultra-large datasets. This paper discusses the optimization techniques for the various operators in DLRM and which component of the systems are stressed by these different operators. The presented techniques are applicable to a broader set of DL workloads that pose the same scaling challenges/characteristics as DLRM.

## I. INTRODUCTION AND BACKGROUND

Over the last two years, there has been steep increase in the number of Bird of a Feather (BoFs) sessions or workshops on high performance deep learning at major supercomputing venues, such as Supercomputing (SC) and International Supercomuputing Conference (ISC). A similar trend can be seen in technical paper tracks focusing on Machine Learning at similar venues. While the presented research is focusing on the extreme-scale and high performance computing (HPC) aspect, it is often limited to training convolutional neural nets (CNN). In essence, researchers study and explore how ResNet50 can be trained on thousands of (accelerated) nodes in less than a minute [1]–[4]. In contrast to this, the "Super 7", namely Facebook, Google, Microsoft, Amazon, Baidu, Alibaba, and Tencent, have published that the CNNs are only contributing a single digit or very low double digit percentage to their workload mix [5]–[7]. Thus, we believe, the HPC research community needs to shift its focus away from CNN models to the models which have the highest percentage in the Super 7's application mix: 1. recommender systems (RecSys) and 2. language models, e.g. recurrent neural networks/long short term memory (RNN/LSTM), and attention/transformer. For the second category, some research has been already published on how to scale it to large platforms [8], [9], whereas RecSys's detailed analyses are basically non-existent in the supercomputing context. This is mainly due to the fact that the industry/the Super 7 have not aligned so far on a standard benchmark (such as ResNet50 for CNNs) for these types of neural networks. Therefore the current release of MLPerf Training [10] is only using a small NCF model. These models do not captures real-life behavior correctly, in terms of model-design and size of the model.

To address these concerns, Facebook recently proposed a deep learning recommendation model (DLRM) [11]. Its purpose is to allow hardware vendors and cloud service providers to study different system configurations, or in simple words to allow for a systematic hardware-software co-design exploration for deep learning systems DLRM comprises of the following major components: a) a sparse embedding realized by tables (databases) of various sizes, and b) small dense multi-layer perceptron (MLP). Both a) and b) interact and feed into c) a larger and deeper MLP. Note that all three parts can be configured (number of features, mini-batch sizes and table sizes) to re-balance and study different topology configurations in a straightforward manner. This work unveils that DLRMs mark the start of a new era of deep learning workloads. In contrast to CNNs, RNNs or Transformers they stress all properties and components of a computing system. This is due to the sparse model portion and as well the dense model portion and of course the communication and interaction between these two. While the sparse portion challenges the memory capacity and bandwidth side, and the dense portion the compute capabilities of the system, the interaction stresses the interconnect. We therefore recognize that a system which maximizes DLRM training performance needs a balanced design between memory capacity, memory bandwidth, interconnect bandwidth and compute/floating point performance.

Facebook open-sourced a simple single process reference implementation in PyTorch and Caffe2 [12]. This implementation yields good performance numbers on GPUs (due to PyTorch GPU-affinity), but CPUs are still a 2nd class citizen in PyTorch. Due to DLRMs' large memory footprint for scale-out scenarios (e.g. real-life terabytes of datasets), we focus on CPUs in this work to understand the scaling implications. On GPUs we would have to scale immediately to very high GPU counts (potentially even spanning embedding tables across multiple GPUs) due to the very limited capacity of

HBM/GDDR per GPU. The problem sizes we will tackle, will not fit onto Nvidia DGX boxes in most of the cases. Additionally, at least up to now, it has been challenging to converge DLRM with FP16 using the default optimizers, such as stochastic gradient descent. This is due to the limited accuracy (not enough mantissa bits) of FP16 [13]. However, a novel scheme that utilizes the BFLOAT16 datatype is introduced in this work. It demonstrates convergence to state-of-the-art accuracies while using DLRM's default optimizer. Therefore, one of the major selling point of GPU's for deep learning, the FP16 tensor cores, cannot be easily leveraged in DLRM with its default optimizers. A comparison between CPU and GPU using FP32 will be covered for a smaller dataset on a single-socket/single GPU basis.

Our contributions are:

- We present a bottom-up performance analysis of standalone micro-apps to identify even more optimization/performance potential in the DLRM code if being freed from DL framework restrictions, and to inform future PyTorch/framework enhancements
- We present a top-down performance analysis and PyTorch code changes that speed-up the reference DLRM topology by roughly two orders of magnitude (110×) on single socket latest Intel Xeon processors. The changes and improvements are upstreamed to PyTorch as part of an active pull-request on github at the time of this writing
- In the DLRM code, we enable hybrid parallelism (running the sparse embedding workload portion in a model-parallel and the MLPs in data-parallel fashion) for multi-process parallelization, targeting multi-node Xeon experiments with optimized communication libraries
- We discuss and analyze the scaling properties of our code on a fat shared-memory system (octo-socket Intel Xeon with 224 cores, to emulate Facebook's recently announced Zion training platform [6]) and traditional HPC clusters (dual-socket Intel Xeon with 56 cores and HPC interconnect, Intel Omnipath (OPA)). We present strong and weak-scaling experiments for various datasets, including the benchmark proposed to MLPerf [10].
- Last but not least, we introduce a novel implementation of the stochastic gradient (SGD) optimizer targeting mix-precision training, called Split-SGD-BF16. It avoids the need for master weights in case of BFLOAT16 mixed precision training. Master weights are nearly always needed in case of FP16 training which complicates the software development. This new SGD future-proofs our work toward next-generation CPUs and accelerators, and is transferable to all other deep learning topologies such as CNNs, RNN/LSTMs, etc.

These contributions are discussed in the following structure: We start with summarizing the DLRM architecture in Sect. II. In terms of covering the mentioned optimization, we split our work into single-socket, multi-socket implementation and performance analysis (Sect. III) through (Sect. VI). Before concluding our research, we will introduce the novel
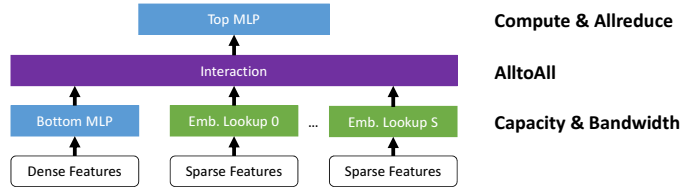


Fig. 1: Schematic of the DLRM topology. It comprises of MLPs and Embedding table look-ups and the corresponding interaction operations. Thus, it stress all important aspects of the underlying hardware platform at the same time: compute capabilities, network bandwidth, memory capacity and memory bandwidth. This is rather unusual for classic HPC applications.

Split-SDG-BF16 as an outlook to next-generation CPUs in Sect. VII. Finally in Sect. VIII and Sect. IX we conclude our work.

## II. THE DLRM BENCHMARK

In this section, we briefly summarize the DLRM benchmark as proposed and published by Facebook [6], [11]. As already touched upon in the introduction, the DLRM benchmark aims at capturing Facebook's most important deep learning training workload through a customize-able benchmark that can be shared with the research community outside of Facebook.

Figure 1 depicts a bird's eye-view of the application. DLRM combines sparse and dense features, which are normally provide through the network in case of online inference, to provide a recommendation for these inputs, e.g. to recommend a specific product for a specific user. The dense features are piped through a classic neural network which is in this case a model problem of a multi-layer perceptron identified as Bottom MLP. MLPs are a very classic concept of neural nets and when using them in a minibatched fashion they transfer computationally into a sequence of matrix multiplications (GEMM) followed by an activation function calculation such as ReLU or Sigmoid. As the activation function calculation is element-wise, it is complexity-wise irrelevant. This can be also seen when examining runtime statistics. As the Bottom MLP is normally relatively small it is somewhat compute bound due to its GEMM-nature but also demands relatively high memory bandwidth as the GEMMs are small and have little re-use. In terms of multi-socket parallelization, MLPs are parallelized best using data-parallelism. This strategy leads to splitting the minibatch across different sockets and requires an allreduce of the weight-gradients during backward propagation, cf. [1]–[4].

The sparse features, which handle categorical data, are mapped through embeddings into a dense space such that they can be combined with the output of the Bottom MLP. We can think of embedding as mulit-hot encoded look-ups into an embedding table $W \in \mathcal{R}^{MxE}$, with $M$ denoting the entries and $E$ the length of each entry. That means, the multi-hot weight-vector $\alpha^T = [0, \ldots, a_{p_1}, \ldots, a_{p_k}, \ldots, 0]$ with elements $a_p \neq 0$ for $p = p_1, \ldots, p_k$ and 0 elsewhere with $p$ being the index for the corresponding lookup items. When adding the minibatch $N$, we can rewrite the embedding

look-up as $L = A^T W$ with $L \in \mathcal{R}^{NxE}$ and the sparse matrix $A = [a_1, \ldots, a_N]$. For more details we refer to [14]. We further want to note that full rows are read from $W$, thus frameworks provide special dense vector routines for embedding lookups and they do not rely on sparse matrix multiplication for its implementation. As shown in Fig. 1, there are normally many dense embedding tables and they can have millions of entries. Therefore the combined task's performance is dictated by available memory capacity and bandwidth.

The results of the Bottom MLP and the embedding look-ups are combined in a so-called Interaction which produces a condensed output. Concat is an example of simple interaction operation. More commonly a self dot product is used as an interaction op which translates to a batched matrix-matrix multiplication as a key kernel. As we can see this operation is parallelized best over the different embedding tables. Therefore, a personalized all-to-all communication is needed within the interaction operation to switch from model parallel to data parallel execution. Finally, the condensed data is fed into a wider and deeper MLP, the Top MLP. Most of the dense floating operations of the execution are spent here since it is merely a GEMM. As in case of the Bottom MLP, this MLP is parallelized best using data-parallel execution.

For computing the loss, DLRM utilizes the classic cross-entropy loss function which is not further discussed as it does not result into any performance implications at all. Therefore, it is not even depicted in Fig. 1.

In terms of computation, we can summarize that the following patterns require attention including concrete performance expectations:

- **GEMM:** the sizes are often (especially in case of scaling out) non-squared, so we need the best possible GEMM routines, and data layouts in order to reach the compute-bound limit of the machine.
- **Embedding Look-ups:** this is a GUPS-like [15] kernel, however we read several consecutive cachelines from memory and sum them up. The expectation is that these operations run at close to peak bandwidth performance of the machine.
- **Allreduce:** for reducing the weigth gradients in the backward pass of the MLPs we need to use the best allreduce algorithm and ensure overlap with the GEMM compute.
- **Alltoall:** for switching between data and model parallelism during the interaction operation

All other kernels and primitives such as activation functions and loss computations are not covered in this work in detail as they do not take a huge fraction of the time and can be relatively easily overlapped or fused with other operations, e.g. ReLU can directly happen inside a custom GEMM routine when the C matrix is still hot in caches [16].

This description points out clearly that DLRMs stress all components of the (cluster) system: compute and network, memory bandwidth and network, capacity and compute. Typical HPC applications as well as CNN training usually stress only one or two aspects. More specifically, DLRM feels like

---

**Algorithm 1** Forward pass of EmbeddingBag Layer

**Inputs**: Weight $W[M][E]$, Indices $I[N+1]$, Offsets $O[NS]$
**Outputs**: Output $Y[N][E]$
1:  **for** $n = 0 \ldots N$ **do**
2:    $start = O[n]$
3:    $end = O[n+1]$
4:    $Y[n][:] = 0.0$
5:    **for** $s = start \ldots end$ **do**
6:      $ind = I[s]$
7:      **for** $e = 0 \ldots E$ **do**
8:        $Y[n][e] \mathrel{+}= W[ind][e]$

---

**Algorithm 2** Backward pass of EmbeddingBag Layer

**Inputs**: GradOut $dY[N][E]$, Offsets $O[N+1]$
**Outputs**: GradWeight $dW[NS][E]$
1:  **for** $n = 0 \ldots N$ **do**
2:    $start = O[n]$
3:    $end = O[n+1]$
4:    **for** $s = start \ldots end$ **do**
5:      **for** $e = 0 \ldots E$ **do**
6:        $dW[s][e] \mathrel{+}= dY[n][e]$

---

running HPCG [17] and HPL [18] in a single and tightly coupled benchmark. This mix of computational motifs poses yet another challenge in designing balanced AI-oriented systems and impacts the machine of choice. Therefore, we will also later study unconventional hardware platforms such as 8 socket servers.

## III. IMPLEMENTATION AND OPTIMIZATIONS - SINGLE SOCKET

In this section we present characteristics of key compute kernels in DLRM, potential issues with their efficient implementation and our bottom up approach for evaluating different alternatives. We then try to integrate these back into PyTorch framework wherever possible to evaluate end-to-end performance gains. In few cases, we just show the estimated benefits of our optimizations if made free from framework restrictions.

### A. Bottom-Up - Embedding Look up Layer

In a typical DL framework, a training iteration consist of 3 passes. A forward pass which computes loss with respect to model parameters, a backward pass which computes gradients with respect to model parameters and an update pass (typically run by an optimizer) which takes learning rate and computed gradients from backward pass and updates the model parameters. Algorithms 1, 2 and 3 shows the forward, backward and the update pass for Embedding table look up and sparse update operation in typical DLRM code when using a linear solver like SGD. The forward and backward pass can be trivially parallelized using `#pragma omp parallel for` over the first line of both forward and backward pass. Compiler can easily vectorize these operations over E loop. The update operation is also trivially simple. However, parallelizing it over NS loop introduces a potential race condition due to indirect access on output rows. Essentially, we need atomicity while performing accumulation to avoid the race condition. Intel Xeon processors provide atomic operations only on integer data types. However, we need floating point atomic_add. There are multiple approaches to solve this problem as discussed below:

**Algorithm 3** Update pass of Sparse EmbeddingBag Layer

**Inputs**: GradWeight $dW[NS][E]$, Indices $I[NS]$m $\alpha$
**InOut**: Weight $W[M][E]$
1: **for** $i = 0 \ldots NS$ **do**
2:     $ind = I[i]$
3:     **for** $e = 0 \ldots E$ **do**
4:        Perform Atomic
5:        $W[ind][e] \mathrel{+}= \alpha * dW[i][e]$

**Algorithm 4** Race Condition Free Update for Sparse EmbeddingBag Layer

**Inputs**: GradWeight $dW[NS][E]$, Indices $I[NS]$m $\alpha$
**InOut**: Weight $W[M][E]$
1: **for** each thread $tid$ in $nThreads$
2:     $M\_start = (M * tid)/nThreads$
3:     $M\_end = (M * (tid + 1))/nThreads$
4:     **for** $i = 0 \ldots NS$ **do**
5:        $ind = I[i]$
6:        **if** $ind >= M\_start$ and $ind < M\_end$ **then**
7:           **for** $e = 0 \ldots E$ **do**
8:              $W[ind][e] \mathrel{+}= \alpha * dW[i][e]$

1) Floating point atomic_add using atomic xchg
2) Using Intel Restricted Transactional Memory (RTM)

Even though we prefer RTM approach over atomic xchg as it allows to use SIMD vector multiply-add instructions inside critical section, we expect both of these approaches perform similarly due to their memory bandwidth limited nature when most of the indices are unique and there is little contention. However, since both of these approaches require a cache line to be modified in its own cache, if there are significant repeated use of indices, there is a potential issue due to excessive cache line thrashing across the core caches when multiple entries of same index are distributed across cores. To address this issue, we implemented the update pass using a race free algorithm and expect it to perform better in a contentious environment/setup. The idea is shown in Algorithm 4. We simply divide the total rows in table among available threads, and each thread goes over the full list of indices but updates the rows only if it belongs to its own range of rows. This algorithm not only eliminates the race conditions but also helps to improve locality of accesses. However, this approach has a potential load imbalance issue if indices are clustered rather than being distributed uniformly across all the rows.

Note that in theory the backward and update pass for EmbeddingBag can be fused into a single operation yielding better performance. However, such fusing are not currently implemented by deep learning framework principles. Nevertheless, we did such an implementation in standalone that achieves up to 1.6x speed-up for embedding updates. However, due to difficulty of integration, we will not consider it any further in this work.

*B. Bottom-Up - Multilayer Perceptron (MLP)*

A MLP consists of (at least three) *fully connected* layers of neurons: the topology starts with an input layer, followed by a number of hidden layers which conclude to the output layer. For the rest of this section we consider the optimization of the *fully connected* layers since they constitute the cornerstone of MLP. Mathematically, an input layer $x \in \mathbb{R}^C$ is mapped to an output layer $y \in \mathbb{R}^K$ via the relation $y = W \cdot x$, where $W \in \mathbb{R}^{K \times C}$ is the weight tensor of the connections between

**Algorithm 5** Forward pass of Fully Connected Layer

**Inputs**: Weight $W[K_b][C_b][b_c][b_k]$, Input $X[C_b][N_b][b_n][b_c]$
**Outputs**: Output $Y[K_b][N_b][b_n][b_k]$
1: Based on $thread\_id$ calculate $K_b\_start$, $K_b\_end$, $N_b\_start$ and $N_b\_end$ to assign output work items
2: **for** $ib_n = N_b\_start \ldots N_b\_end$ **do**
3:     **for** $ib_k = K_b\_start \ldots K_b\_end$ **do**
4:        /* Prepare batch-reduce GEMM arguments */
5:        **for** $ib_c = 0 \ldots C_b - 1$ **do**
6:           $A_{ptrs}[ib_c] = \&W[ib_k][ib_c][0][0]$
7:           $B_{ptrs}[ib_c] = \&X[ib_c][ib_n][0][0]$
8:        $Out = \&Y[ib_k][ib_n][0][0]$
9:        **batchreduce_gemm**$(A_{ptrs}, B_{ptrs}, Out, C_b)$

the neurons. During the training process, $N$ multiple inputs ($N$ is the so-called mini-batch size) are grouped together yielding the equation $Y = W \cdot X$ with $W \in \mathbb{R}^{K \times C}$, $X \in \mathbb{R}^{C \times N}$ and $Y \in \mathbb{R}^{K \times N}$. Observe that by increasing the mini-batch $N$, we fundamentally increase the weight tensor reuse. Typical implementations of Fully Connected layers (e.g. the one in the PyTorch DLRM version) leverage a large GEMM call and they apply the activation functions onto the GEMM outputs. Even though such an approach is straightforward to implement, its performance can be underwhelming for two reasons: i) typical high-performance GEMM library calls internally perform packing of sub-matrices to ameliorate TLB misses and cache conflict misses [19], and ii) the multi-threaded implementation of GEMM with shapes arising from small mini-batch values $N$ may not fully exploit the available data reuse.

In this section we dive into the details of the forward propagation algorithm of the MLP training process (also used for inference); we also implemented all the required kernels of the back-propagation training in an analogous fashion. Our MLP implementation which leverages the batch-reduce GEMM microkernel follows the same principles of previous work [20], and we present it here for completeness. However, in this work we modify the tensor layouts since we found this variation to yield better performance in training where the activations tensors in the backward by weights pass are the analogous of the weights tensors in the forward/backward by data passes. We emphasize upfront here that incorporating blocked tensor layouts in frameworks (for performance reasons) is rather cumbersome, yet we examine with standalone MLP training code the potential upside.

Algorithm 5 shows the implementation of the forward propagation in the training process of fully connected layers. First, we highlight the blocked tensor layout; all the 2-dimensional tensors are transformed into 4-dimensional ones by blocking the mini-batch dimension $N$ with a factor $b_n$ and the tensor dimensions $C$ and $K$ with blocking factors $b_c$ and $b_k$ respectively. Such a blocked layout exposes better locality and avoids large, strided sub-tensor accesses which are known to cause TLB misses and cache conflict misses in case the leading dimensions are large powers of 2. Also, in contrast to previous work [20], we use the $[C_b][N_b][b_n][b_c]$ format for activations since it provides similar performance benefits in the backward by weights training pass.

Our algorithm first assigns the output sub-tensor blocks to the available threads (line 1) and every thread then for
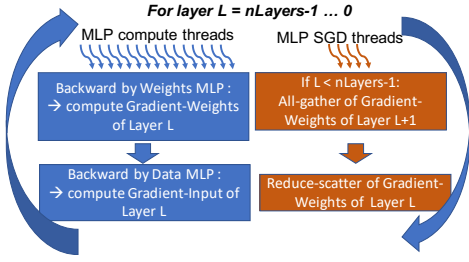
**For layer L = nLayers-1 ... 0**

MLP compute threads | MLP SGD threads

Backward by Weights MLP :
→ compute Gradient-Weights
of Layer L

If L < nLayers-1:
All-gather of Gradient-
Weights of Layer L+1

Backward by Data MLP :
→ compute Gradient-Input of
Layer L

Reduce-scatter of Gradient-
Weights of Layer L

Fig. 2: Overlapping the SGD solver with the MLP GEMMs in the back-propagation pass of MLP.

each assigned output $Y$ block calculates the addresses of the $W$ and $X$ sub-tensor blocks that need to be multiplied and reduced onto the current output $Y$ block (lines 5-7). Note that our batch-reduce GEMM kernel is JIT-ed, and allows small values of blocking values $b_n$ to be used, and as such we can extract parallelism from the mini-batch dimension even for small values of $N$. By following the loop ordering of Algorithm 5, a weight sub-tensor is reused by each thread $N_b\_end - N_b\_start - 1$ times, potentially from some level of cache. Also, multiple threads are able to read weights from shared caches when the assigned $Y$ blocks correspond to the same subspace of the $K$ dimension. Finally, in case a weight sub-tensor does not fit in the targeted/desired level of cache, we can further block loops at lines 3 and 5. These cache blocking techniques in combination with the flexible blocking factors $b_n$, $b_c$ and $b_k$ which yield high performance micro-kernels, attempt to maximize the data reuse in the GEMM. The backward by data and backward by weights passes follow the same Algorithm 5 since they also instantiate GEMM operations.

## IV. IMPLEMENTATION AND OPTIMIZATIONS - MULTI SOCKET

### A. Bottom-Up - Multi-socket and multi-node MLP training

A major challenge during multi-socket and multi-node MLP training is to hide the allreduce communication behind the GEMMs of the backward pass. Therefore we study this problem in a standalone fashion, freed from any framework limitations.

In order to scale the MLP training to multiple sockets and multiple nodes we adopted the following strategy: First we employ data parallelism (i.e. minibatch parallelism) across sockets/nodes, and within each socket we rely on a 2D GEMM decomposition as described in Section III-B. In such a setup, the only communication required across sockets/nodes is during the SGD whereas the 2D decomposition of GEMMs within the socket strives for optimal data movement/reuse.

For the GEMM computations, we leveraged our high-performance Fully Connected/GEMM layers described in subsection III-B. In order to minimize the overhead of the communication in the SGD, we overlapped the SGD solver with the back-propagation MLP kernels. More specifically, the SGD solver can be seen as an All-reduce operation among all the Weight Gradient tensors residing in all sockets/nodes. In

our implementation we materialize the all-reduce operation via a reduce-scatter and an all-gather operation. Figure 2 visualizes how we can overlap the all-gather and the reduce-scatter operations with the compute MLP kernels. Given $T$ cores/threads per socket, we dedicate $S$ threads per socket for the SGD/communication of the gradient weights, and $T - S$ threads for the computation/GEMMs in the back-propagation of the Fully Connected Layers. We tune the value of $S$ in order to balance the communication time in SGD and the computation time in GEMMs. Such MPI progression threads are well-known in large-scale HPC applications [21]–[24].

### B. Top-Down - All-to-all Communication in PyTorch Multi-Process

The original DLRM code from Facebook has support for device-based hybrid parallelization strategies. It follows data-parallelism for MLP and model-parallelism for embeddings. This means that MLP layers are replicated on each agent and the input to first MLP layer is distributed over minibatch dimension whereas embeddings tables are distributed across available agents and each table produces output worth full minibatch. This leads to mismatch in minibatch size at interaction operation and requires communication to align minibatch of embedding output with bottom MLP output. The multi-device (e.g. socket, GPU, accelerator, etc.) implementation of DLRM uses sequence of scatters, one per table, to distribute embedding output over minibatch before entering to interaction operation. We extend this parallelization strategy to multi-process using MPI, where a rank can be considered equivalent to a device. We use the Distributed Data Parallel (DDP) module to wrap bottom and top MLPs, whereas for embeddings we simply distribute tables across available ranks. We call this approach "ScatterList". One of its drawbacks is that, it makes multiple calls to the communication backend (one call per table). When there are multiple tables per rank this approach turns out to be inefficient, and we optimize it by coalescing output of multiple local tables into one buffer and invoking just one scatter per rank. We call this approach as "Fused Scatter". However, if we look at it carefully, this is a well-known all-to-all communication pattern known to HPC world for years. Nevertheless, DL frameworks such as PyTorch used to lack primitives for supporting this communication pattern. PyTorch has recently added experimental support for alltoall primitive to their distributed backend. We used it in our distributed DLRM code to further minimize the number of calls to communication backend to make it just one call. We call this as "Alltoall" version of our code.

### C. Top-Down - Optimal Communication Backend in PyTorch

To get good scaling efficiency we need to overlap communication with computation. To enable asynchronous communication, the MPI backend of PyTorch spawns a separate thread to drive the communication. The idea is to enable our Bottom-Up approach directly in PyTorch to demonstrate its value in end-to-end runs. The master thread simply enqueues the requested operation to a MPI thread and waits for completion when it is
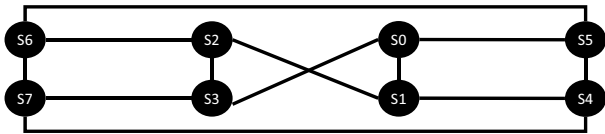
Fig. 3: Schematic of an 8 socket Intel Xeon system, an Inspur TS860M5 machine. As each socket offers only 3 UPI links, the sockets are arranged in a twisted hypercube in order to offer a balanced communication path between all sockets.
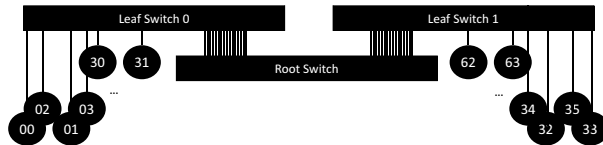


Fig. 4: Schematic of the 64 socket cluster which is connected through a pruned fat-tree. 32 sockets connect with full bandwidth to leaf switches and these two switches are connected with half bandwidth through a root switch.

ready to consume the output of communication. While using a separate thread is good for driving communication, as we discussed in previous section, we need multiple threads to saturate the full communication bandwidth. Also, it is very important to manage CPU affinities for this thread. Otherwise it can interfere with compute threads and potentially cause performance degradation. We found that Intel's oneCCL [25] library tries to solve these issues. However, it is not integrated into PyTorch as built-in communication backend. PyTorch has recently added an experimental feature that supports adding custom communication backend for torch.distributed module. OneCCL is integrated into PyTorch using this feature [26]. We call it as CCL backend and our Alltoall version as "CCL-Alltoall" when running with the CCL backend. Finally, for debugging communication performance, we modified DDP module to optionally perform blocking allreduce and added autograd profiling hooks to the communication backends. This allows us to measure time spent in communication primitives which is important for detailed performance analysis.

## V. HARDWARE PLATFORMS AND BENCHMARKING SETUP

As motivated in the introduction we will focus on CPU-based platforms as the sparse embedding requires often a huge memory footprint. Due to DLRM's huge memory consumption, a single standard dual-socket high-performance computing node is only useful for some small problems and tuning DLRM's kernels. Therefore we decided to use two different test beds. The first one is normally used by big database companies such as Oracle or SAP, and is aligned with Facebook's recently announced Zion platform [6]: an eight socket Skylake Xeon platform. The second one is a classic HPC configuration: 32 dual-socket Cascade Lake Xeon nodes connected by an Intel Omnipath (OPA) pruned fat-tree with two rails per node.

### A. 8 Socket Shared-Memory Node

The simplest way of running an application with a large memory footprint is to use node with a lot of DRAM. Normally, these nodes are used for databases as the accesses happen randomly and cannot be easily located.

Figure 3 depicts the layout of an Intel Xeon Skylake/Cascade Lake system featuring 8 sockets. The Platinum series processor offers 3 point-to-point Ultra Path Interconnect (UPI) links. However, each socket needs to communicate with 7 neighbors. Therefore, the sockets are organized in a twisted hypercube fabric. That ensures that 3 neighbors can be reached

in one hop and the remaining 4 neighbors in two hops. Each of the UPI link offers roughly 22 GB/s bidirectional bandwidth which is comparable with a 100G fabric in cluster installations, resulting into an aggregated system UPI bandwidth of 260 GB/s as the machine has 12 unique UPI connections. Each socket is exposed as a NUMA node in the system such that software can easily optimize the various distances and latencies in case of different local and remote NUMA memory accesses.

In our case, we are using an Inspur TS860M5 machine which has a 2x4 socket node design. Each socket is equipped with an top-of-the line Intel Xeon Skylake processor (SKX), the Intel Xeon Platinum 8180. It features 28 cores at an AVX512 turbo frequency of 2.3 GHz and 1.7 AVX512 base frequency. Therefore each socket provides 4.1 TFLOPS FP32 peak performance. With respect to memory, the system is optimized for high bandwidth and high capacity by using 12 dual-rank 16 GB DDR4-2400 DIMMs per socket offering 100 GB/s memory bandwidth. In total the machine offers 224 core providing 32 FP32-TFLOPS at 800 GB/s bandwidth with a capacity of 1.5 TB.

### B. 64 Socket Cluster Architecture

We are going to compare the performance of the 8 socket machine to a traditional HPC installation, namely several dual-node systems connected by low-latency and high-bandwidth fabric, in our case Intel OPA.

The layout of the cluster we leverage for this experiment, is depicted in Fig. 4. Each of the 64 sockets houses its own OPA network adapter which offers 100G connectivity at 1us latency. As OPA switches allow to connect only 48 endpoints we implemented a classic pruned fat-tree: 16 nodes with 32 sockets are connected to one switch each, and then both leaf switches are connected with 16 links to the root switch. Storage and head nodes are also connected to the root switch. This results into 200 GB/s within each 32 sockets' leaf and 200 GB/s between the leaves as we prune the bandwidth with a ratio of 2:1 when going up to the root.

In terms of CPU per socket, the cluster is very similar to the 8 socket system covered earlier. Each socket houses a top-bin Intel Cascade Lake processoer (CLX), the Intel Xeon Platinum 8280 which offers 100Mhz additional clock over the 8180 of the 8 socket system. Its 28 cores run at an AVX512 turbo frequency of 2.4 GHz and 1.8 AVX512 base frequency. Therefore each socket provides 4.3 TFLOPS FP32 peak performance. With respect to memory, the system

| Configuration Parameter | Small | Large | MLPerf |
|---|---|---|---|
| Minibatch (MB) (Single Socket) ($N$) | 2048 | - | 2048 |
| Global MB for Strong Scaling ($GN$) | 8192 | 16384 | 16384 |
| Local MB for Weak Scaling ($LN$) | 1024 | 512 | 2048 |
| Avg. look-ups per Table ($P$) | 50 | 100 | 1 |
| Number of Tables ($S$) | 8 | 64 | 26 |
| Embedding Dimension ($E$) | 64 | 256 | 128 |
| #rows per table ($M$) | $1 \cdot 10^6$ | $6 \cdot 10^6$ | Up to 40M |
| Length Inputs Bottom MLP | 512 | 2048 | 13 |
| #Layers Bottom MLP | 2 | 8 | 3 |
| Bottom MLP Size | 512 | 2048 | 512-256-128 |
| #Layers Top MLP | 4 | 16 | 4 |
| Bottom MLP Size | 1024 | 4096 | 512-512-256-1 |

TABLE I: DLRM model specifications used in this work. Small is taken from the release paper of DLRM [11].

is optimized for high bandwidth only as its main target is HPC. Therefore it only uses 6 dual-rank 16 GB DDR4-2666 DIMMs per socket offering 105 GB/s memory bandwidth. In total the machine offers 1,792 cores providing 275 FP32-TFLOPS at 6.7 TB/s bandwidth with a capacity of 6 TB. 4 of the 32 nodes have 192GB/socket memory allowing us large single socket runs.

## C. Difference and Similarities Between Both Systems

As we can see from the two descriptions, we have designed and chosen two very different multi-socket systems in terms of technology. However, the application visible performance specifications are very similar. While OPA offers 100G between sockets, data still needs to be copied through the network card stack which means multiple internal data copies. In contrast UPI offers as well 100G connectivity but we can copy data without any additional movements. Therefore on the 8 socket platform, we can use non-temporal (non-cached) write flows of full cachelines to minimize the communication volume to a bare minimum. These flows can be regarded as true one-side communication.

## D. Benchmarking Setup

In order to evaluate different pressure points of our hardware systems (covered in the following section) we will use three different configurations of DLRM in this work summarized in Tab. I. The Small variant is identical to the model problem used in DLRM's release paper [11]. Large variant is the small problem scaled in every aspect for scale-out runs. The MLPerf configuration is recently proposed as a benchmark config for performance evaluation of a recommendation system training [27]. It uses Criteo Terabyte Dataset [28] for its convergence run. For single socket run, we use `numactl` to bind it to socket 0 whereas distributed run always use 1 socket per rank and occupy the node first before going multiple nodes.

## VI. ANALYSIS AND RESULTS

In this Section we present the performance results for the aforementioned implementations. We start with a bottom-up analysis and end with large-scale cluster experiments. Each result is the average of multiple batches of iterations to account for potential small run-to-run performance variations.
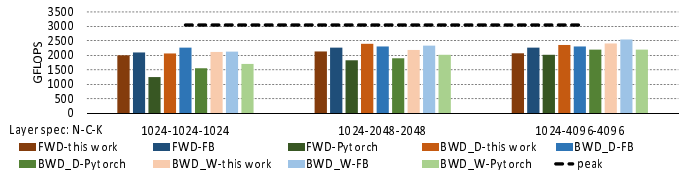


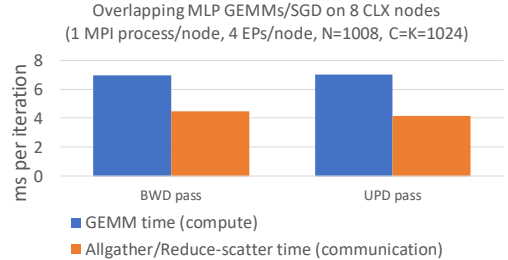Fig. 5: MLP training kernel performance on single SKX socket



Fig. 6: Communication/Computation overlap on 8 CLX nodes (1 MPI process/node, N=1008, C=K=1024)

## A. Bottom-Up: Single Socket standalone MLP results

Figure 5 shows the performance of 5-layer MLP training kernels on a single SKX socket. In particular we compare all three training passes, namely forward (FWD), backward by data (BWD_D) and backward by weights (BWD_W) for three different implementations: i) this work (orange-shaded bars), ii) an MLP training code optimized by Facebook [29] for the SKX platform (blue-shaded bars) and iii) the PyTorch implementation that leverages multi-threaded MKL calls for the involved GEMMs (green-shaded bars). Also, we experimented with various MLP configurations, where we fixed the minibatch $N$ to be 1024 and we varied the feature map dimensions $C$ and $K$ (we tried the values 1024, 2048 and 4096). The Facebook MLP implementation is similar to ours: it employs a NUMA-aware and thread-aware blocking of the Fully Connected layers/GEMMs and for the single-threaded kernels it uses serial MKL GEMM calls. We observe that our approach and Facebook's implementation show similar performance on a single socket: the average performance across all configurations and all passes is 72% and 75% of peak respectively. On the contrary, the MLP implementation in PyTorch which employs large, multi-threaded MKL calls shows average efficiency 61% of peak and is $\sim 18\%$ slower than our approach.

## B. Bottom-Up: Multi Socket / Multi-node standalone MLP results

We evaluate the efficacy of the computation/communication overlapping during the standalone MLP training implementation. For the 8 socket shared memory node, we implemented the SGD via non-temporal flows, and the communication takes place over the UPI connections between the sockets. More specifically, we dedicated 4 cores per socket for the SGD communication flow which is overlapped with the back-propagation GEMMs that leverage 24 cores per socket. Similarly, for the multinode setup, we assign one MPI-rank per socket, dedicated 24 cores per socket for the back-
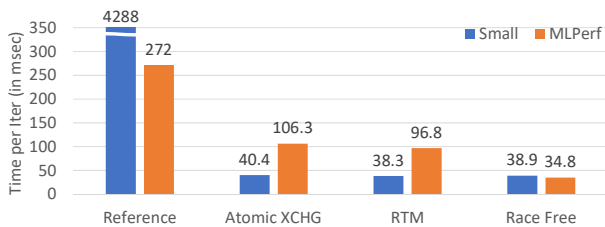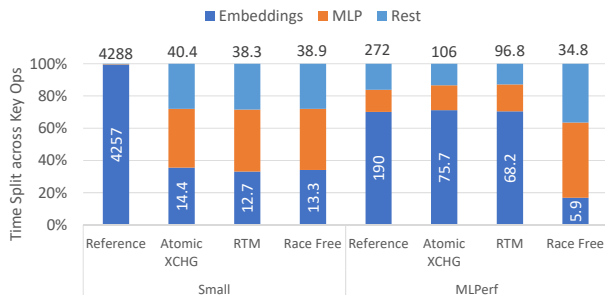
Fig. 7: DLRM single socket performance



Fig. 8: DLRM single socket performance breakdown

propagation GEMMs whereas we created 4 MPI Endpoints (EPs) per socket for the communication. As we see on Figure 6, such a setup/problem configuration is sufficient to completely hide the communication (orange bars) behind the computations/GEMMs (blue bars). It is also worth-noting that our overlapping scheme can hide the SGD/communication cost over UPI (i.e. in the 8 socket node). E.g. for the smallest config of Figure 5 the backward by data and backward by weights GEMMs take 5.40 and 5.39 ms respectively while the corresponding communication operations that are overlapped require 2.84 and 1.86 ms.

## C. End-to-End: Single Socket DLRM results

We started our DLRM analysis using the latest release of PyTorch v1.4.0. After analyzing profiling data for each operation of the DLRM code, we found that 99% of the execution time was spent in just one kernel related to the sparse embedding look up. This was caused by a naive CPU backend implementation which was focused on functionality instead of performance. We found few more similar instances but performance impact was dominated by a single embedding update kernel. Therefore, we carefully reviewed all PyTorch framework operations involved when executing DLRM. While trivial issues were fixed in PyTorch source code, we implemented a PyTorch c++ extension for EmbeddingBag operations for quick experimentation and implemented various update strategies described in Sec. III. The performance benefits depend on config and ranges from $110\times$ for small config to $8\times$ for the MLPerf config. We could not run the large config on single socket as it needs minimum of 450GB DRAM memory capacity. Figure 8 shows percent of time spent in various key components before and after applying optimizations. We can see that for the small config, all the 3 embedding update strategies perform more or less same as indices distribution is random and have very little contention. However, for the

| Parameter | Small | Large | MLPerf |
|---|---|---|---|
| Mem capacity required for all tables (in GB) | 2 | 384 | 98 |
| Minimum sockets required | 1 | 4 | 1* |
| Maximum Ranks to scale | 8 | 64 | 26 |
| Total AllReduce Size (in MB) | 9.5 | 1047 | 9.0 |
| Strong Scaling Alltoall Volume (in MB) | 15.8 | 1024 | 208 |

TABLE II: DLRM model characteristics for distributed run

MLPerf config, we observe a lot of contention with the terabyte dataset causing up to $10\times$ slowdown for embedding update operation compared to our race free algorithm. We also observe that even though it appears (with a first glance) that DLRM performance is embedding-heavy, after optimizations, it takes about 30% of total time for small config matching it with MLP time whereas for the MLPerf config, embeddings take less than 20% of total time.

Last but not least, we can summarize that the small config, runtime per iteration decreased from 4288ms to 38ms. Additionally, we want to mention that previous work [11] timed an NVIDIA V100 GPU at 62ms running exactly the same DLRM problem using the Caffe2 frontend. However, to be fair, we have to mention that V100 has roughly $3.5\times$ more FP32-FLOPS than Skylake/Cascade and $8\times$ more available bandwidth at much smaller memory capacity. Therefore, we can expect a fully-optimized GPU software stack to be at around 10-15ms for the small problem, being 2-3x faster than our optimized single-socket CPU version. Nevertheless, that puts the CPU into a very strong position as it can run virtually any problem configuration without memory limitations as we will show in the next section. More importantly, a CPU-cluster is not limited by small memory capacity per socket and therefore the sweet spot between MLP and embedding performance can be easily picked.

## D. End-to-End: Multi Socket / Multi-node DLRM results

As a part of distributed DLRM scaling analysis, we wanted to find answers to following questions and reason about the performance we see.

1) What is the pure compute scaling?
2) What is the cost of pure communication when there is no compute interference?
3) How does the allreduce and alltoall communication time change as we increase the number of ranks?
4) What is the impact of overlapping communication?
5) How much communication time can be overlapped?
6) How much time is spent in preparing actual communication such as copy to flat buffers or gradient averaging?

We evaluated the multi-node/multi-socket performance of DLRM for all the 3 configs on our 64 socket cluster and an 8 socket shared-memory node. We performed strong scaling experiments to size reduction in time-to-train as we leverage more sockets to solve a fixed problem. We also performed weak scaling analysis to understand the pure communication overheads as we scale on multiple sockets. Due to pure model parallelism used for distributing embedding tables, the maximum number of ranks we can use to scale a config depends on the total number of embedding tables we have in the config.

Since, the small config has only 8 embedding tables, we can scale it on up to 8 sockets. Large config can be scaled on up to 64 sockets while MLPerf can be scaled on up to 26 sockets. Moreover, since the large config has a large memory foot print, we can only start running it with minimum of 4 sockets. Therefore, we use 4 ranks best performance (CCL-Alltoall version) as our baseline for computing efficiency or speed up for the large config. For the other two configs, we use optimized single socket performance as a baseline for comparison. Before looking at scaling performance we want to understand the communication characteristics of the 3 configs and set some expectations.

$$SZ_{\text{allreduce}} = \sum_{l=0}^{n} f_i^l \times f_o^l + f_o^l \qquad (1)$$

$$SZ_{\text{alltoall}} = S \times N \times E \qquad (2)$$

Eq. 1 shows the size for allreduce as seen by each rank where $f_i^l$ and $f_o^l$ are the input and output feature maps for given MLP layer and $n$ is total number of layers in top and bottom MLP. $SZ_{\text{allreduce}}$ is independent of number of ranks or minibatch size. Therefore, cost of allreduce increases steadily as we increase the number of ranks and imposes a major challenge for strong scaling compared to weak scaling. Eq. 2 shows total volume of alltoall communication, now across all the ranks, related to embedding tables communication. This volume is proportional to the global minibatch ($N$) and thus remains constant for strong scaling and increases proportional to number of ranks for weak scaling. Since, this volume gets divided across all the ranks, the size of point-to-point message reduces $4\times$ as we double the number of ranks. Also, alltoall directly benefits multiple links such as UPI across sockets and OPA across nodes. Similarly, cost of alltoall is expected to reduce by 4x when going from 2 to 4 ranks and then steadily go down as we increase number of ranks further. Table II shows the allreduce buffer size for MLPs and alltoall volume for the strong scaling configs we used in our evaluation. Looking at these parameters, we expect the small and large problem to be allreduce-bound whereas the MLPerf config would initially be alltoall-bound and becomes allreduce-bound for high rank counts. Finally, since the output of allreduce is used at the end of backward pass, it can be overlapped over entire compute of backward pass (see Section IV-A) whereas alltoall can only be overlapped with bottom MLP compute. Thus, cost of alltoall is more difficult to hide compared to allreduce cost. In order to answer the questions 1 to 6 effectively, we instrumented the PyTorch source code to optionally perform blocking communication for allreduce and alltoall (shown in the "blocking" part of the graphs in Figure 10, 11, 13 & 14). The bars in "overlapping" part shows compute and exposed communication time as seen by application after overlap.

*1) Strong Scaling:* Figure 9 shows strong-scaling speed-up and efficiency of scaling when scaled up to 64 ranks (64R). We achieve up to 8.5x end-to-end speed up for the MLPerf config when running on 26 sockets (33% efficiency) and about 5x-6x speed up when increasing the number of sockets by

8x for the small and large configs ($\sim$60%-71% efficiency). It is obvious that the native alltoall performs better than scatter-based alltoall implementations and yields more than 2x performance benefits. However, even MPI-based alltoall leaves a lot of performance on the table as it becomes evident from performance of CCL-Alltoall version which gives up to 1.4x additional speed up for end-to-end time. To understand this in more detail, we look at the compute-communication breakdown for large and MLPerf configs as shown in Figure 10. Here we skip the analysis of small config which has a similar behavior to that of the large problem. As we can see, when using the MPI-backend, not only communication time is higher than for the CCL-backend but even compute times grow significantly. This is not intuitive. After looking into individual kernels' performance, we observed that almost all compute kernels were slowed down due to communication overlap. This was happening because the thread spawned by PyTorch MPI backend for driving MPI communication was interfering with compute threads and slowing down both the compute and communication. To avoid this, CCL provides a mechanism to bind the communication threads to specific cores and we exclude these cores from compute cores. Therefore, we do not see compute slowdown issues but much better overlap of compute and communication with the CCL-backend. Figure 11 depicts further breakdown of communication cost as time spent in pre- and post-processing for actual communication and actual time spent in wait calls. The pre- and post-processing costs remain comparable across both backends. However, even the pure cost of communication is lower with CCL-backend compared to the MPI-backend. This is because CCL uses multiple cores to drive the communication. Another puzzling thing we observed for large problem is the cost of alltoall. We see a huge alltoall cost for MPI backend when overlapping communication is used but almost negligible cost when blocking communications are used. After careful analysis, we concluded that this is happening due to in-order completion nature of MPI-backend that shows up as cost of allreduce at alltoall wait.

Given all this, we see good initial efficiency of scaling for small and large configs which drops steadily due to exposed allreduce cost as we increase the number of ranks. For the MLPerf config, the profile starts with lower efficiency at 2 ranks due to very high alltoall cost. Then initially, efficiency goes up until 8 ranks as cost of alltoall reduces and again drops at 16 and 26 ranks as allreduce cost starts to assemble the majority of communication. This is in line with our expectations for our simple performance with Eq. 1 and Eq. 2.

*2) Weak Scaling:* Figure 12 shows weak-scaling speed-up and efficiency. With weak-scaling, we achieve up to 17x end-to-end speed up for the MLPerf config when running on 26 sockets (65% efficiency) and 13.5x speed up (84% efficiency) when increasing number of sockets by 16x (64 ranks compared to 4 ranks optimized baseline) for the large config. In case of the small config, we achieve about 6.4x speed up on 8 sockets (80% efficiency). We see a similar trend for using native alltoall compared to scatter based alltoall implementations
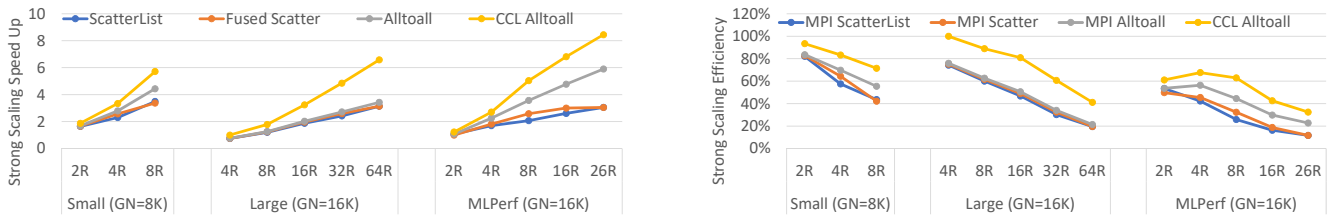
Fig. 9: DLRM strong scaling performance wrt. optimized baseline (Left: speed-up, Right: scaling efficiency)
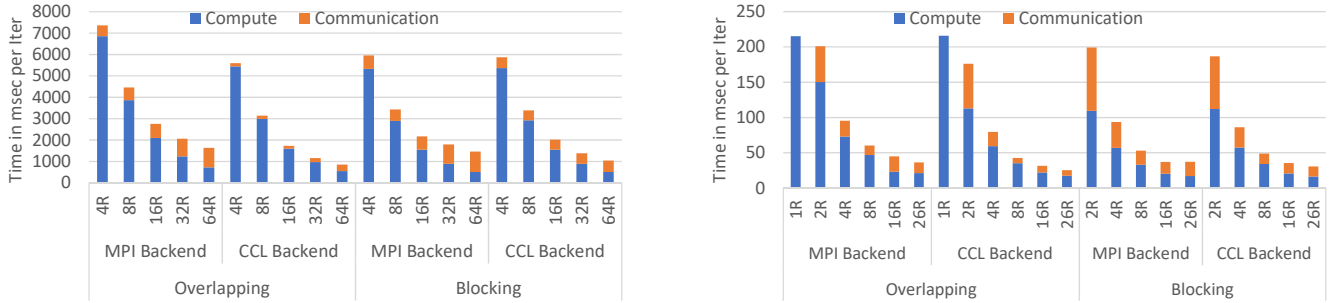


Fig. 10: Compute-Communication time break up with and without overlap for strong scaling (Left: Large Config, Right: MLPerf Config)
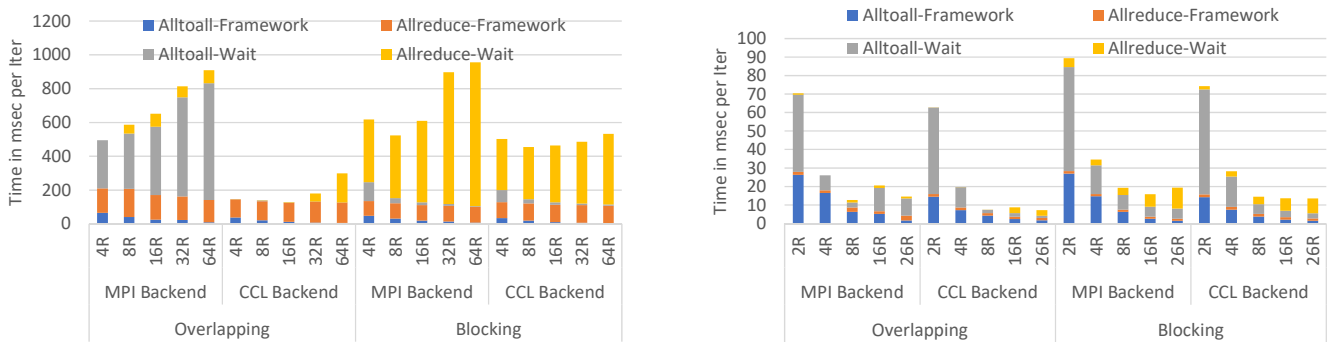


Fig. 11: Communication time break up with and without overlap for strong scaling (Left: Large Config, Right: MLPerf Config)
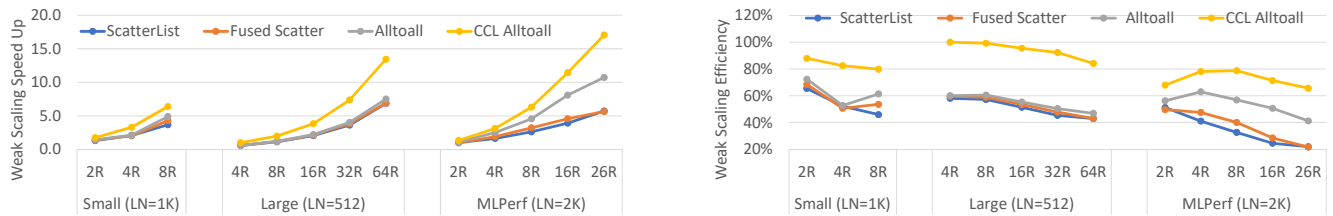


Fig. 12: DLRM weak scaling performance wrt. optimized baseline (Left: speed-up, Right: scaling efficiency)

and MPI-based alltoall vs. CCL-Alltoall version as has been observed for strong scaling. E.g. we continue to see increased compute time due to communication overlap for MPI backend as shown in Figure 13. For the CCL backend the difference in compute time with and without overlap is almost negligible. Also, for the MLPerf config, we see, how cost of communication goes down at first (up to 8 ranks) and then increases again due to different characteristics of alltoall and allreduce.

This can be further confirmed from a cost break-down for alltoall and allreduce depicted in Figure 14. Finally, we see the compute cost (even with blocking communication) for the MLPerf config increases slowly as we weak-scale, instead we expect it to remain constant. A more detailed look revealed that this additional cost comes from the current data loader design which always reads the data for full global minibatch on each rank and with weak scaling that cost steadily grows.
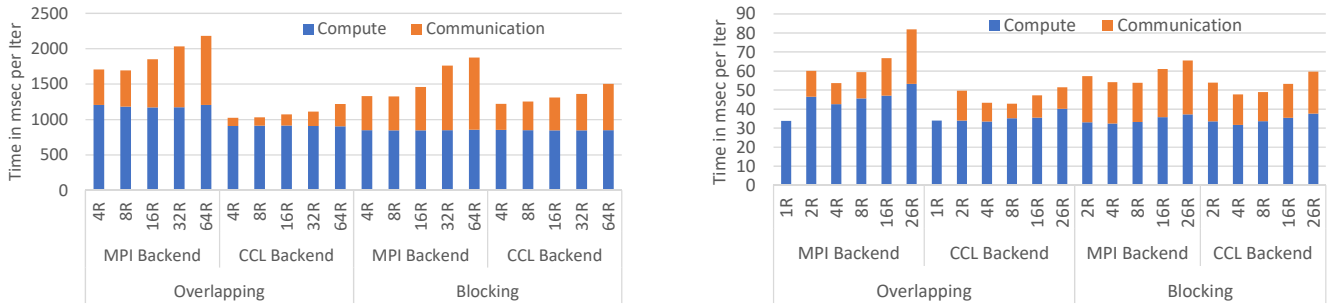
Fig. 13: Compute-Communication time break up with and without overlap for weak scaling (Left: Large Config, Right: MLPerf Config)
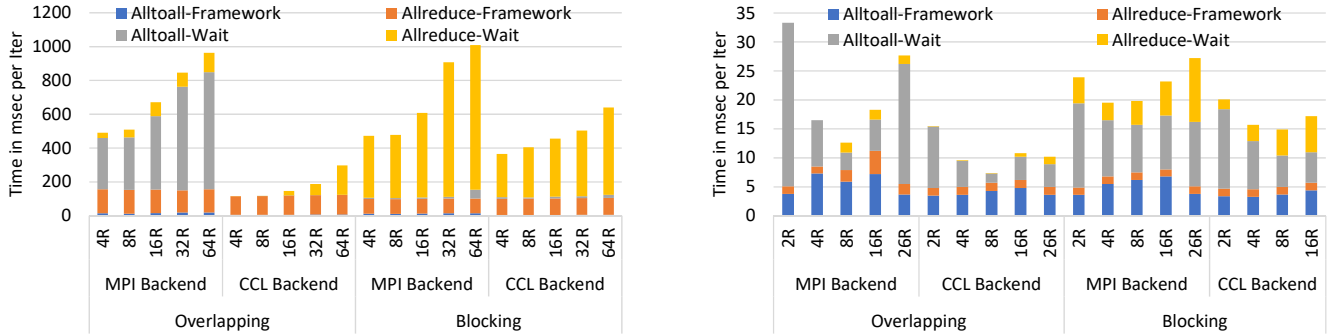


Fig. 14: Communication time break up with and without overlap for weak scaling (Left: Large Config, Right: MLPerf Config)

For small and large configs, we use random dataset which does not account for time spent in data loader.

*3) Scaling on 8-Socket Shared Memory System:* Finally, in Fig. 15, we discuss the strong scaling performance of the 8 socket shared-memory system which exhibits similar behavior to a cluster with similar node count. The only difference is, that the cost of alltoall does not decrease from 4 to 8 sockets as expected. This is because the alltoall implementation is not optimally tuned for twisted-hypercube connectivity on the system and so links are not utilized optimally. Additionally, even optimal algorithms would need multiple rounds of communication such that only 1.5x can be expected when comparing 4 to 8 nodes. This is evident for the MLperf config. Nevertheless, as the 8 socket system fits all our workloads under investigation, it can be seen as a small cluster in an appliance form-factor as no external high-performance fabric is needed. So we can understand why Facebook prefers such a design in a loosely-coupled, Ethernet-based cloud datacenter [30], [31].

## VII. SPLIT-SGD FOR NATIVE BFLOAT16 TRAINING

BFLOAT16 is a new, but no-standard, floating-point format [32], [33] that is gaining traction due to its ability to work well in machine learning algorithms, in particular deep learning training. In contrast to the IEEE754-standardized 16bit (FP16) variant, BFLOAT16 does not compromise at all on range when being compared to FP32. As a reminder, FP32 numbers have 8 bits of exponent and 24 bits of mantissa (one implicit). BFLOAT16 cuts 16 bits from the 24-bit FP32 mantissa to create a 16-bit floating point datatype. In contrast FP16, roughly halves the FP32 mantissa to 10 explicit bits and has to reduce the exponent to 5 bits to fit the 16-bit datatype envelope. BFLOAT16 therefore perfectly aliases with the upper half of IEEE754-FP32 numbers.

Classic training approaches using BFLOAT16, FP16 or even int16 require so-called master-weights, a full precision FP32 copy of the entire model throughout the full training process [34]–[36]. With 16bit "regular" weights, this results into a 200% (or 3X) overhead in storing the model. As the model size in case of DLRM is mainly determined by the embedding tables, and we are already starved for capacity, a mixed precision solver with FP16 would require hundreds of Gigabytes more capacity than a vanilla FP32 approach.

Instead, our Split-SGD-BF16 solver aims at efficiently exploiting the aforementioned aliasing of BF16 and FP32. Therefore it reduces the overhead and being equal to FP32 training, wrt. to capacity requirements, as master weights are implicitly stored. The trick is that we do not store FP32 values as a single tensor in a classic fashion. Instead, we split them into their high and low 16bit parts. First we store all 16 MSBs of the FP32 numbers and then all 16 LSBs of the numbers as two separate tensors. The 16 MSBs represent a valid BFLOAT16 number and we store it as part of model. We exclusively use those in the forward and backward passes, whereas lower bits are only required in optimizer so we store them as additional state in optimizer and the actual update uses both, MSBs and LSBs and runs therefore a fully FP32-accurate
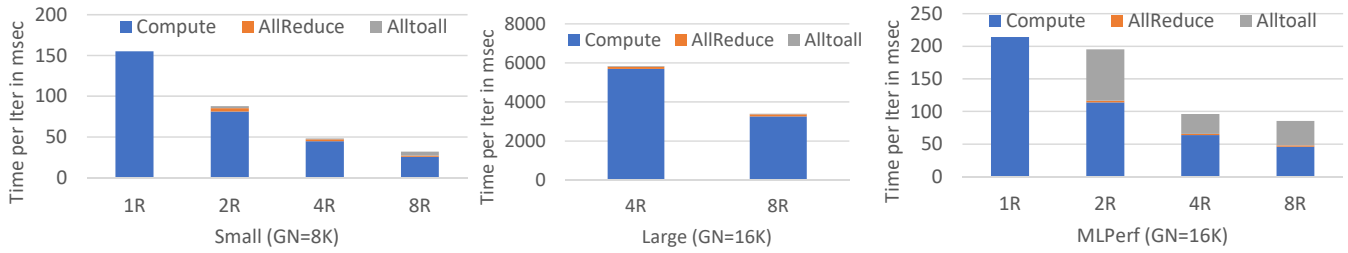
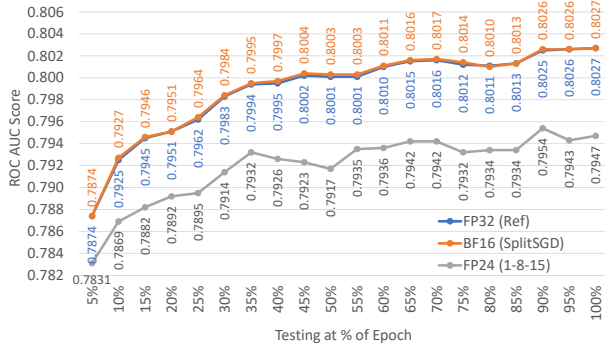Fig. 15: Strong scaling performance on 8-socket shared-memory system



Fig. 16: Training accuracy with Mix precision BF16

update. That means 66% of the training passes enjoy a 2x bandwidth reduction without negatively affecting the accuracy. Figure 16 demonstrates that this technique allows us to train the DLRM MLPerf configuration to state-of-the-art with an error of less than 0.001%. Please note, in this plot we also use a bit-accurate software emulation of Intel Xeon Cooper Lake BFLOAT16 dot-product instructions `vdpbf16ps` [37]. When silicon is available, this will help to also significantly speed-up the MLP portions as well. The described concept of Split-SGD-BF16 is independent of the workload and can be transferred to any BFLOAT16 deep learning training task.

In addition, we tried to run with only 8 additional LSBs. However, these are not enough to train DLRM to required accuracy. Note, that we tried to replicate FP16 embedding training with stochastic quantization as described in [13]. Unfortunately, we were not able to train various DLRM configuration to state-of-the-art using SGD. We believe, more work is needed to generalize [13] beyond the very simple model problems presented.

## VIII. RELATED WORK

While we believe this paper is the first of its kind in covering large-scale DLRM training with an in depth analysis, some high level considerations for DLRM training were published in [38]. Apart from training, several studies have analyzed inference properties of DLRM [39], [40] and this even pertains to FPGA acceleration [41]. All these works have in common, that the CPU code, used as a baseline, is not as highly optimized as ours. With respect to large-scale training, results for CNNs have been shown here [1]–[4] and RNN/LSTMs are covered in [8], [9]. Large-scale HPC and Deep Learning

on on-premises systems and in the cloud accelerated by MPI progressions threads has been previously discussed in [21]–[24]. The efficient implementation of various deep learning compute operations on Intel CPUs such as convolutions and LSTM cells were discussed in [8], [16], [20] and our work is based on these findings. BFLOAT16 training with master weights on CPUs and without Split-SGD has been shown for a variety of workloads in [34]. For completeness we have to note that NVIDIA also open-sourced a recommender system architecture which seems to be more aligned with NCF and less complex than DRLM [42].

## IX. CONCLUSION AND FUTURE WORK

In this paper we presented how programming methodologies and hardware architectures from the field of HPC can be used to significantly speed-up training AI topologies, specifically the Deep Learning Recommender Models (DLRM). On Sky-lake/Cascade Lake CPUs we achieve a performance boost of more than two orders of magnitude (110×) on a single socket and studied strong and weak scaling with excellent results for various problems sizes. This is true for large shared-memory nodes and clusters with up to 64 CPU sockets. In addition to making the DLRM benchmark framework fit for CPU cluster architectures of today, we also introduced and enabled DLRM with a novel SGD optimizer. This leverages the BFLOAT16 datatype which is soon to be supported by various CPU architectures while matching FP32 accuracies. Last but not least we demonstrated that a single socket CPU is 2x faster than a single V100 GPU, while we expect the GPU could be in theory 2x-3x faster if the code would be fully optimized. As the CPU cluster is not limited by capacity per socket, the best sweet spot between MLP and embedding performance, or alltoall and allreduce, respectively, can be picked.

## REFERENCES

[1] N. Dryden, N. Maruyama, T. Moon, T. Benson, M. Snir, and B. Van Essen, "Channel and filter parallelism for large-scale cnn training," in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, ser. SC '19. New York, NY, USA: ACM, 2019, pp. 10:1–10:20. [Online]. Available: http://doi.acm.org/10.1145/3295500.3356207

[2] T. Akiba, S. Suzuki, and K. Fukuda, "Extremely large minibatch SGD: training resnet-50 on imagenet in 15 minutes," *CoRR*, vol. abs/1711.04325, 2017. [Online]. Available: http://arxiv.org/abs/1711.04325

[3] Y. You, Z. Zhang, C. Hsieh, and J. Demmel, "100-epoch imagenet training with alexnet in 24 minutes," *CoRR*, vol. abs/1709.05011, 2017. [Online]. Available: http://arxiv.org/abs/1709.05011

[4] H. Mikami, H. Suganuma, P. U.-Chupala, Y. Tanaka, and Y. Kageyama, "Imagenet/resnet-50 training in 224 seconds," *CoRR*, vol. abs/1811.05233, 2018. [Online]. Available: http://arxiv.org/abs/1811.05233

[5] N. P. Jouppi, C. Young, N. Patil, D. A. Patterson, G. Agrawal, R. Bajwa, S. Bates, S. Bhatia, N. Boden, A. Borchers, R. Boyle, P. Cantin, C. Chao, C. Clark, J. Coriell, M. Daley, M. Dau, J. Dean, B. Gelb, T. V. Ghaemmaghami, R. Gottipati, W. Gulland, R. Hagmann, R. C. Ho, D. Hogberg, J. Hu, R. Hundt, D. Hurt, J. Ibarz, A. Jaffey, A. Jaworski, A. Kaplan, H. Khaitan, A. Koch, N. Kumar, S. Lacy, J. Laudon, J. Law, D. Le, C. Leary, Z. Liu, K. Lucke, A. Lundin, G. MacKean, A. Maggiore, M. Mahony, K. Miller, R. Nagarajan, R. Narayanaswami, R. Ni, K. Nix, T. Norrie, M. Omernick, N. Penukonda, A. Phelps, J. Ross, A. Salek, E. Samadiani, C. Severn, G. Sizikov, M. Snelham, J. Souter, D. Steinberg, A. Swing, M. Tan, G. Thorson, B. Tian, H. Toma, E. Tuttle, V. Vasudevan, R. Walter, W. Wang, E. Wilcox, and D. H. Yoon, "In-datacenter performance analysis of a tensor processing unit," *CoRR*, vol. abs/1704.04760, 2017. [Online]. Available: http://arxiv.org/abs/1704.04760

[6] M. Smelyanskiy, "Zion: Facebook next-generation large memory training platform," in *Proceedings of Hot Chips 31*, ser. Hot Chips 31, 2019.

[7] U. Gupta, X. Wang, M. Naumov, C. Wu, B. Reagen, D. Brooks, B. Cottel, K. M. Hazelwood, B. Jia, H. S. Lee, A. Malevich, D. Mudigere, M. Smelyanskiy, L. Xiong, and X. Zhang, "The architectural implications of facebook's dnn-based personalized recommendation," *CoRR*, vol. abs/1906.03109, 2019. [Online]. Available: https://arxiv.org/abs/1906.03109

[8] D. D. Kalamkar, K. Banerjee, S. Srinivasan, S. Sridharan, E. Georganas, M. E. Smorkalov, C. Xu, and A. Heinecke, "Training google neural machine translation on an intel cpu cluster," in *2019 IEEE International Conference on Cluster Computing (CLUSTER)*, Sep. 2019, pp. 1–10.

[9] Y. You, J. Hseu, C. Ying, J. Demmel, K. Keutzer, and C.-J. Hsieh, "Large-batch training for lstm and beyond," in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, ser. SC '19.   New York, NY, USA: ACM, 2019, pp. 9:1–9:16. [Online]. Available: http://doi.acm.org/10.1145/3295500.3356137

[10] P. R. Mattson, C. Cheng, C. Coleman, G. Diamos, P. Micikevicius, D. A. Patterson, H. Tang, G.-Y. Wei, P. Bailis, V. Bittorf, D. M. Brooks, D. Chen, D. Dutta, U. Gupta, K. Hazelwood, A. Hock, X. Huang, B. Jia, D. Kang, D. E. Kanter, N. Kumar, J. Liao, G. Ma, D. Narayanan, T. Oguntebi, G. Pekhimenko, L. Pentecost, V. J. Reddi, T. Robie, T. S. John, C.-J. Wu, L. Xu, C. Young, and M. Zaharia, "Mlperf training benchmark," *ArXiv*, vol. abs/1910.01500, 2019.

[11] M. Naumov, D. Mudigere, H.-J. M. Shi, J. Huang, N. Sundaraman, J. Park, X. Wang, U. Gupta, C.-J. Wu, A. G. Azzolini, D. Dzhulgakov, A. Mallevich, I. Cherniavskii, Y. Lu, R. Krishnamoorthi, A. Yu, V. Y. Kondratenko, S. Pereira, X. Chen, W. Chen, V. Rao, B. Jia, L. Xiong, and M. Smelyanskiy, "Deep learning recommendation model for personalization and recommendation systems," *ArXiv*, vol. abs/1906.00091, 2019.

[12] Facebook, "Dlrm," https://github.com/facebookresearch/dlrm, 2019.

[13] J. Zhang, J. Yang, and H. Yuen, "Training with low-precision embedding tables," -, 2018.

[14] M. Naumov, "On the dimensionality of embeddings for sparse features and data," -, 2019. [Online]. Available: http://arxiv.org/abs/1901.02103

[15] V. Aggarwal, Y. Sabharwal, R. Garg, and P. Heidelberger, "Hpcc randomaccess benchmark for next generation supercomputers," in *2009 IEEE International Symposium on Parallel Distributed Processing*, May 2009, pp. 1–11.

[16] E. Georganas, S. Avancha, K. Banerjee, D. Kalamkar, G. Henry, H. Pabst, and A. Heinecke, "Anatomy of high-performance deep learning convolutions on simd architectures," in *Proceedings of the International Conference for High Performance Computing, Networking, Storage, and Analysis*, ser. SC '18. Piscataway, NJ, USA: IEEE Press, 2018, pp. 66:1–66:12. [Online]. Available: https://doi.org/10.1109/SC.2018.00069

[17] J. Dongarra, M. A. Heroux, and P. Luszczek, "A new metric for ranking high-performance computing systems," *National Science Review*, vol. 3, no. 1, pp. 30–35, 01 2016. [Online]. Available: https://doi.org/10.1093/nsr/nwv084

[18] A. Petitet, R. Whaley, J. Dongarra, and A. Cleary, "Hpl a portable implementation of the high-performance linpack benchmark for distributed-memory computers," -, 01 2008.

[19] K. Goto and R. A. Geijn, "Anatomy of high-performance matrix multiplication," *ACM Transactions on Mathematical Software (TOMS)*, vol. 34, no. 3, p. 12, 2008.

[20] E. Georganas, K. Banerjee, D. Kalamkar, S. Avancha, A. Venkat, M. Anderson, G. Henry, H. Pabst, and A. Heinecke, "Harnessing deep learning via a single building block," *2020 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, 2020.

[21] A. Breuer, A. Heinecke, and M. Bader, "Petascale local time stepping for the ader-dg finite element method," in *2016 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, 2016, pp. 854–863.

[22] K. Vaidyanathan, K. Pamnany, D. D. Kalamkar, A. Heinecke, M. Smelyanskiy, J. Park, D. Kim, A. Shet, G, B. Kaul, B. Joo, and P. Dubey, "Improving communication performance and scalability of native applications on intel xeon phi coprocessor clusters," in *Proceedings of the 2014 IEEE 28th International Parallel and Distributed Processing Symposium*, ser. IPDPS 14.   USA: IEEE Computer Society, 2014, p. 10831092. [Online]. Available: https://doi.org/10.1109/IPDPS.2014.113

[23] S. Sridharan, K. Vaidyanathan, D. D. Kalamkar, D. Das, M. E. Smorkalov, M. Shiryaev, D. Mudigere, N. Mellempudi, S. Avancha, B. Kaul, and P. Dubey, "On scale-out deep learning training for cloud and hpc," *ArXiv*, vol. abs/1801.08030, 2018.

[24] T. Hoefler and A. Lumsdaine, "Message progression in parallel computing - to thread or not to thread?" in *2008 IEEE International Conference on Cluster Computing*, 2008, pp. 213–222.

[25] Intel, "oneccl," https://github.com/oneapi-src/oneCCL, 2019.

[26] ——, "torch-ccl," https://github.com/intel/torch-ccl, 2020.

[27] "Mlperf: Fair and useful benchmarks for measuring training and inference performance of ml hardware, software, and services." [Online]. Available: http://www.mlperf.org/

[28] "Download terabyte click logs," 12/12/2013. [Online]. Available: https://labs.criteo.com/2013/12/download-terabyte-click-logs

[29] "Facebook-optimized mlp on multisocket systems," Accessed on 12/15/2019. [Online]. Available: https://github.com/jspark1105/tbb_test

[30] K. Hazelwood, S. Bird, D. Brooks, S. Chintala, U. Diril, D. Dzhulgakov, M. Fawzy, B. Jia, Y. Jia, A. Kalro, J. Law, K. Lee, J. Lu, P. Noordhuis, M. Smelyanskiy, L. Xiong, and X. Wang, "Applied machine learning at facebook: A datacenter infrastructure perspective," in *2018 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, 2018, pp. 620–629.

[31] J. Park, M. Naumov, P. Basu, S. Deng, A. Kalaiah, D. S. Khudia, J. Law, P. Malani, A. Malevich, N. Satish, J. Pino, M. Schatz, A. Sidorov, V. Sivakumar, A. Tulloch, X. Wang, Y. Wu, H. Yuen, U. Diril, D. Dzhulgakov, K. M. Hazelwood, B. Jia, Y. Jia, L. Qiao, V. Rao, N. Rotem, S. Yoo, and M. Smelyanskiy, "Deep learning inference in facebook data centers: Characterization, performance optimizations and hardware implications," *ArXiv*, vol. abs/1811.09886, 2018.

[32] "Tensorflow development summit," March 30 2018.

[33] *BFLOAT16 Hardware Numerics Definition*.   Santa Clara, USA: Intel Corporation, 2018.

[34] D. D. Kalamkar, D. Mudigere, N. Mellempudi, D. Das, K. Banerjee, S. Avancha, D. T. Vooturi, N. Jammalamadaka, J. Huang, H. Yuen, J. Yang, J. Park, A. Heinecke, E. Georganas, S. Srinivasan, A. Kundu, M. Smelyanskiy, B. Kaul, and P. Dubey, "A study of bfloat16 for deep learning training," *ArXiv*, vol. abs/1905.12322, 2019.

[35] P. Micikevicius, S. Narang, J. Alben, G. F. Diamos, E. Elsen, D. García, B. Ginsburg, M. Houston, O. Kuchaiev, G. Venkatesh, and H. Wu, "Mixed precision training," *ArXiv*, vol. abs/1710.03740, 2018.

[36] D. Das, N. Mellempudi, D. Mudigere, D. D. Kalamkar, S. Avancha, K. Banerjee, S. Sridharan, K. Vaidyanathan, B. Kaul, E. Georganas, A. Heinecke, P. Dubey, J. Corbal, N. Shustrov, R. Dubtsov, E. Fomenko, and V. O. Pirogov, "Mixed precision training of convolutional neural networks using integer operations," *ArXiv*, vol. abs/1802.00930, 2018.

[37] *Intel Architecture Instruction Set Extensions and Future Features Programming Reference*.   Santa Clara, USA: Intel Corporation, 2020.

[38] M. Naumov, J. Kim, D. Mudigere, S. Sridharan, X. Wang, W. Zhao, S. Yilmaz, C. Kim, H. Yuen, M. Ozdal, K. Nair, I. Gao, B.-Y. Su, J. Yang, and M. Smelyanskiy, "Deep learning training in facebook data centers: Design of scale-up and scale-out systems," *ArXiv*, vol. abs/2003.09518, 2020.

[39] U. Gupta, X. Q. Wang, M. Naumov, C.-J. Wu, B. Reagen, D. Brooks, B. Cottel, K. Hazelwood, B. Jia, H.-H. S. Lee, A. Malevich, D. Mudigere, M. Smelyanskiy, L. Xiong, and X. Zhang, "The architectural implications of facebook's dnn-based personalized recommendation," *2020 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, pp. 488–501, 2019.

[40] U. Gupta, S. Hsia, V. Saraph, X. Q. Wang, B. Reagen, G.-Y. Wei, H.-H. S. Lee, D. M. Brooks, and C.-J. Wu, "Deeprecsys: A system for optimizing end-to-end at-scale neural recommendation inference," *ArXiv*, vol. abs/2001.02772, 2020.

[41] R. Hwang, T. Kim, Y. Kwon, and M. Rhu, "Centaur: A chiplet-based, hybrid sparse-dense accelerator for personalized recommendations," *ISCA'20*, 2020.

[42] NVIDIA, "Hugectr," https://github.com/NVIDIA/HugeCTR, 2019.