

Demystify Communication Behavior in Training Deep Learning Recommendation Models

Ching-Hsiang Chu

Research Scientist

Meta

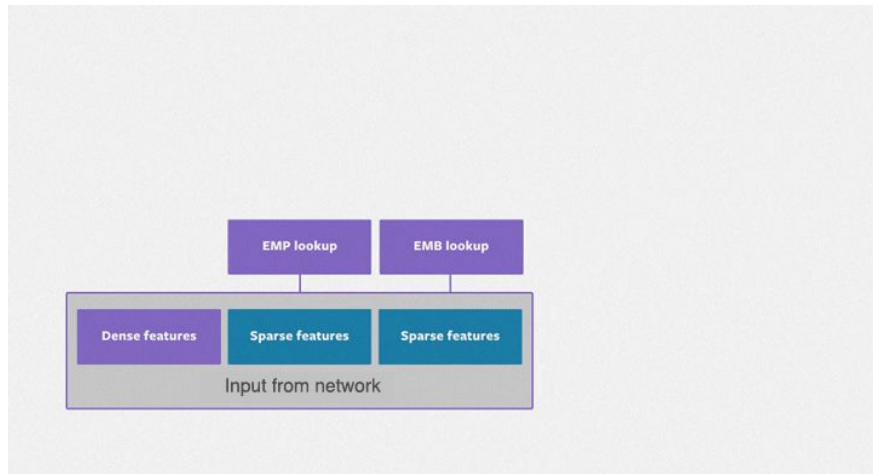
chchu@fb.com

Agenda

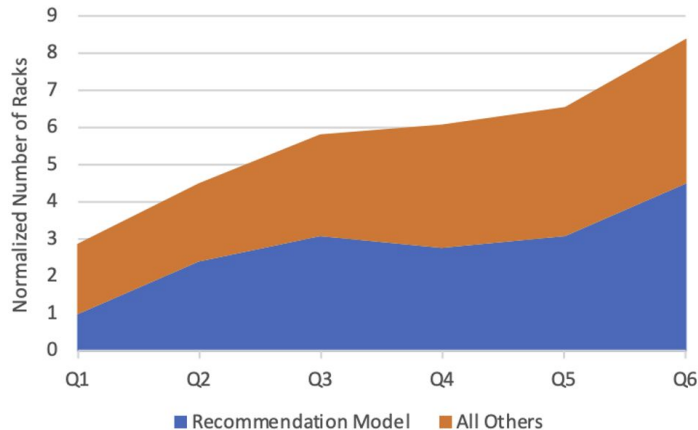
- **Introduction**
- Demystify Communication Behavior in Training DLRM
 - With a real-world example
- Communication Reproduction
- Summary

Deep learning recommendation models (DLRMs)

- DLRMs are used extensively in many companies for building recommendation systems
- MLPerf training and inference benchmarks (<https://mlcommons.org/en/>)

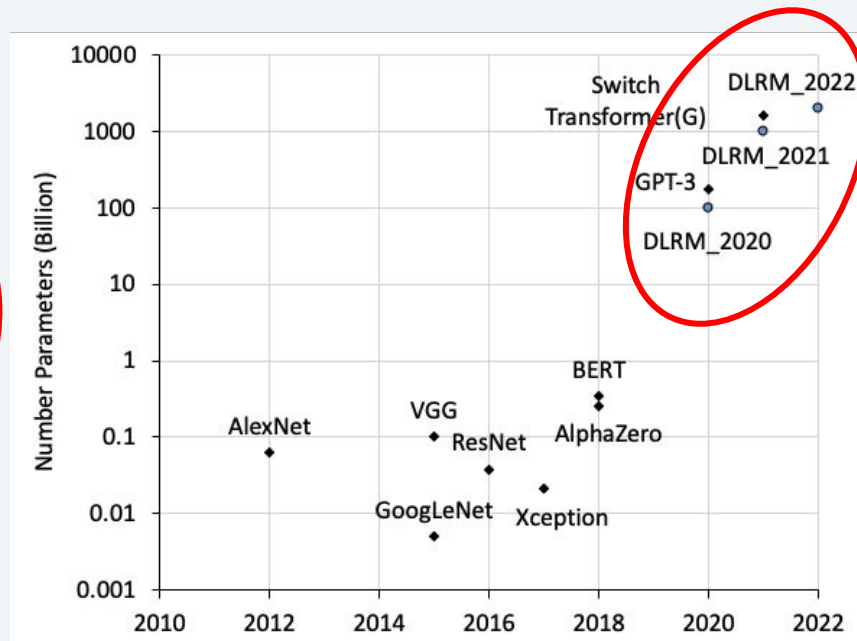
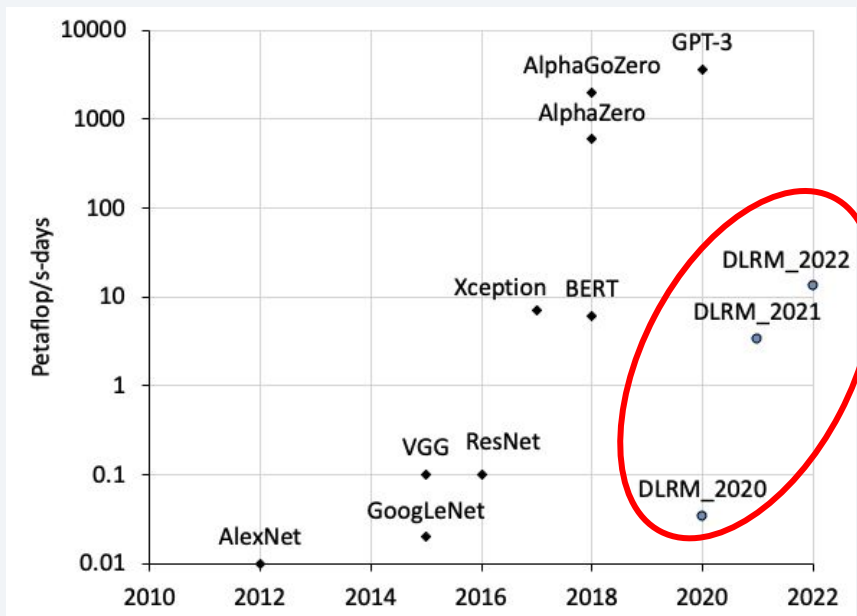


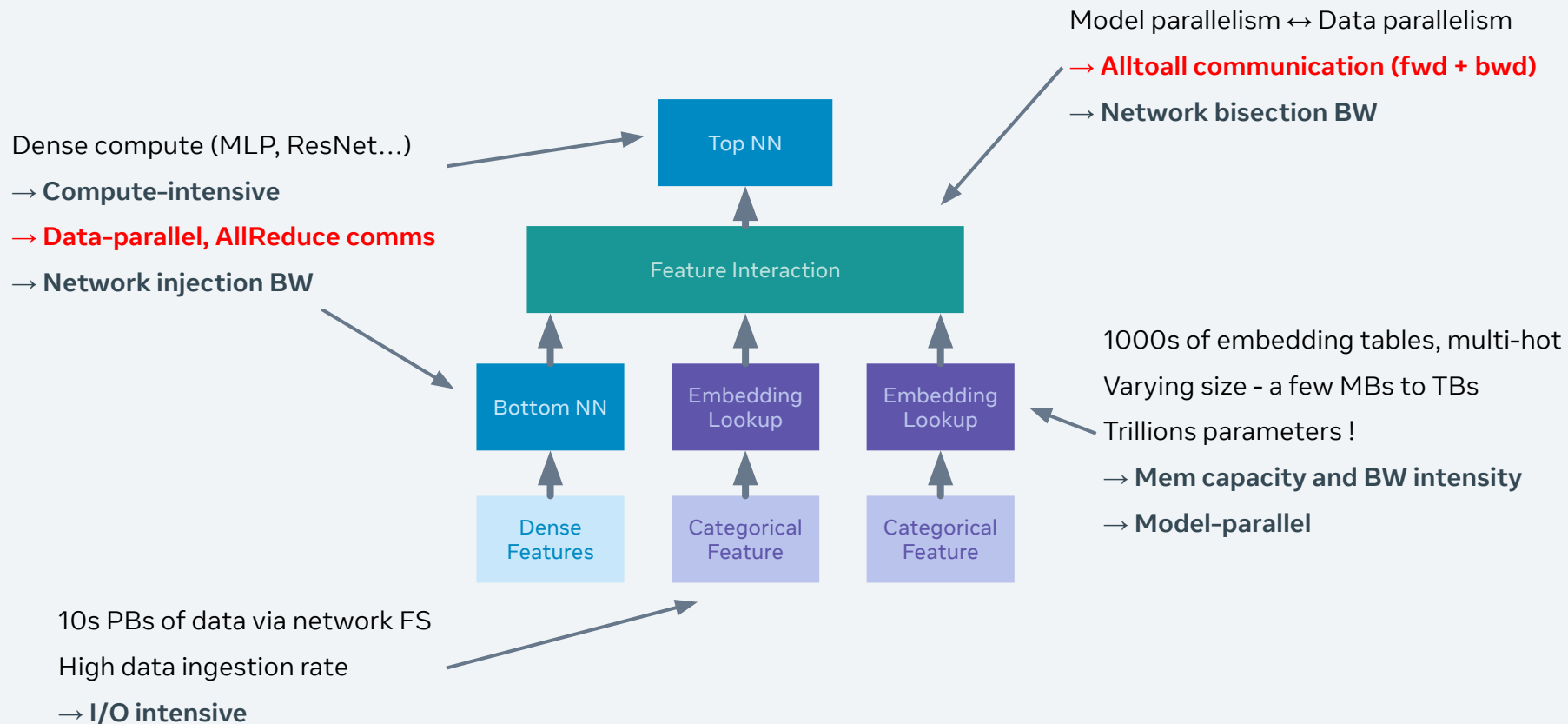
<https://ai.facebook.com/blog/dlrm-an-advanced-open-source-deep-learning-recommendation-model/>



Recommendation models are different !

- Lower compute intensity
- Larger sizes

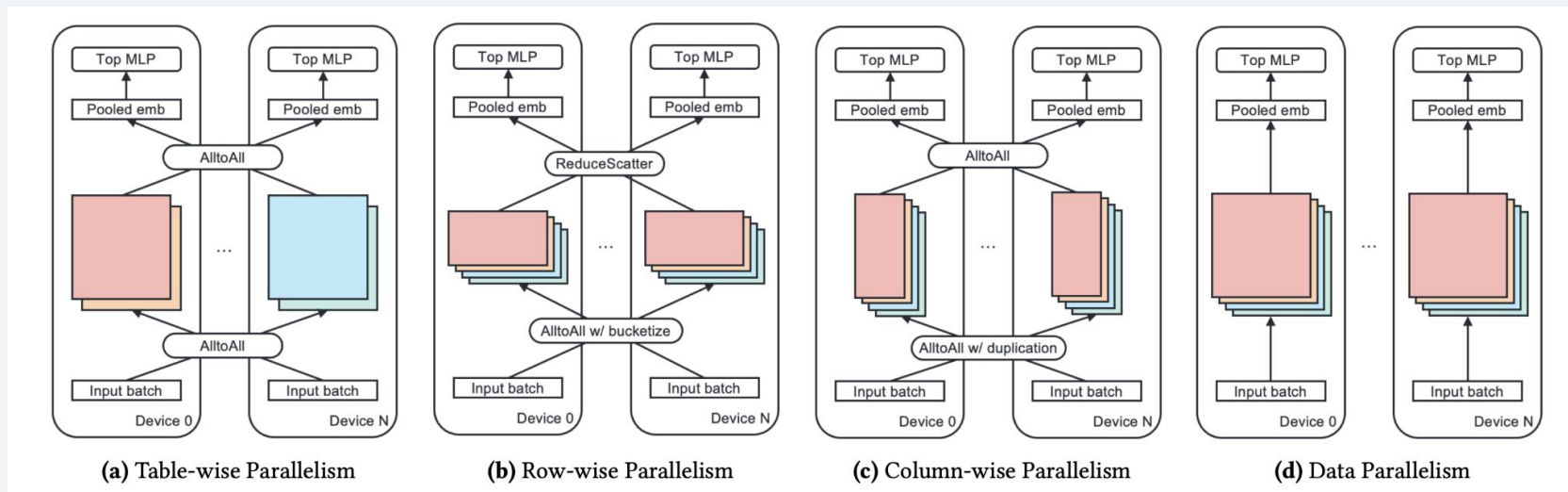




Parallelism - all (possible) ways

Flexible 4D Model parallelism for embedding tables

- Sharding across tables, rows, columns and data
- Hierarchical sharding combining multiple strategies



Agenda

- Introduction
- **Demystify Communication Behavior in Training DLRM**
 - **With a real-world workload**
- Reproduction of Communication
 - Collect Communication trace
 - Replay communication trace
- Summary

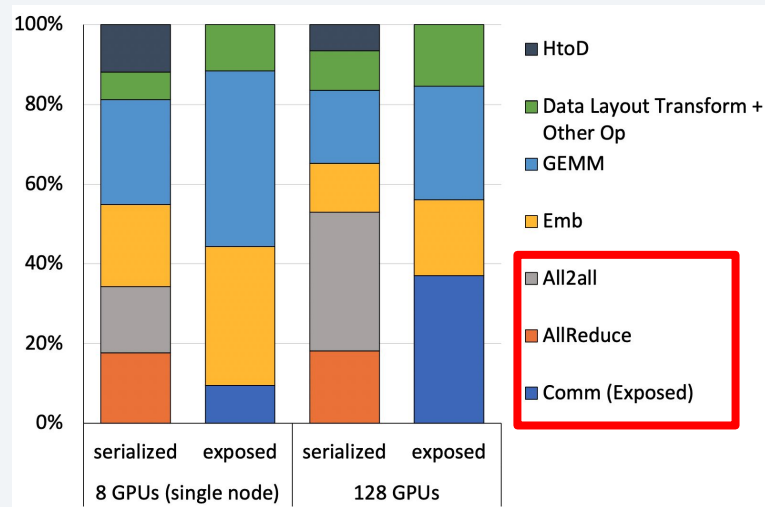
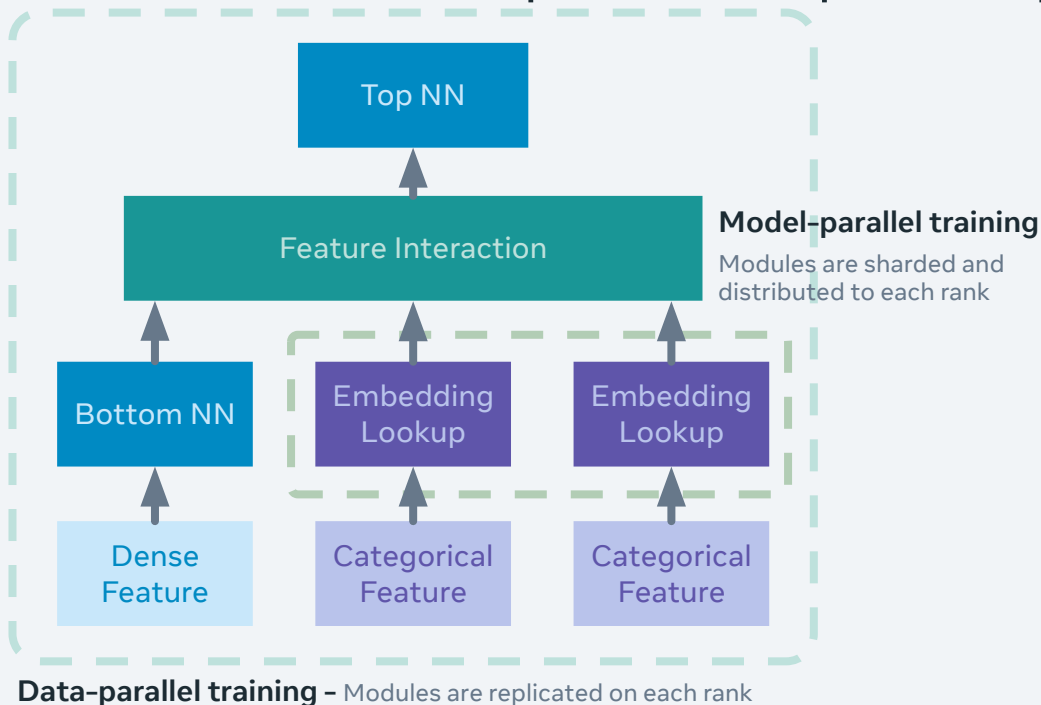
Communication in DLRM-2020 workloads

- Testbed*
 - 128 Nvidia V100 GPUs, 8 GPUs per node
 - 8 *200G RDMA NICs per node
- SW stack
 - DLRM (<https://github.com/facebookresearch/dlrm>)
 - PyTorch (<https://pytorch.org/>)
 - NCCL (<https://github.com/NVIDIA/nccl>)

Model	model-A
Num parameters	793B
MFLOPS per sample	638
Num of emb tables	≈ 1000s
Embedding table dims (range [min, max], avg)	[4, 384] avg: 93
Avg pooling size	15
Num MLP layers	20
Avg MLP size	3375
Target local batch size	512
Achieved QPS	1.2M

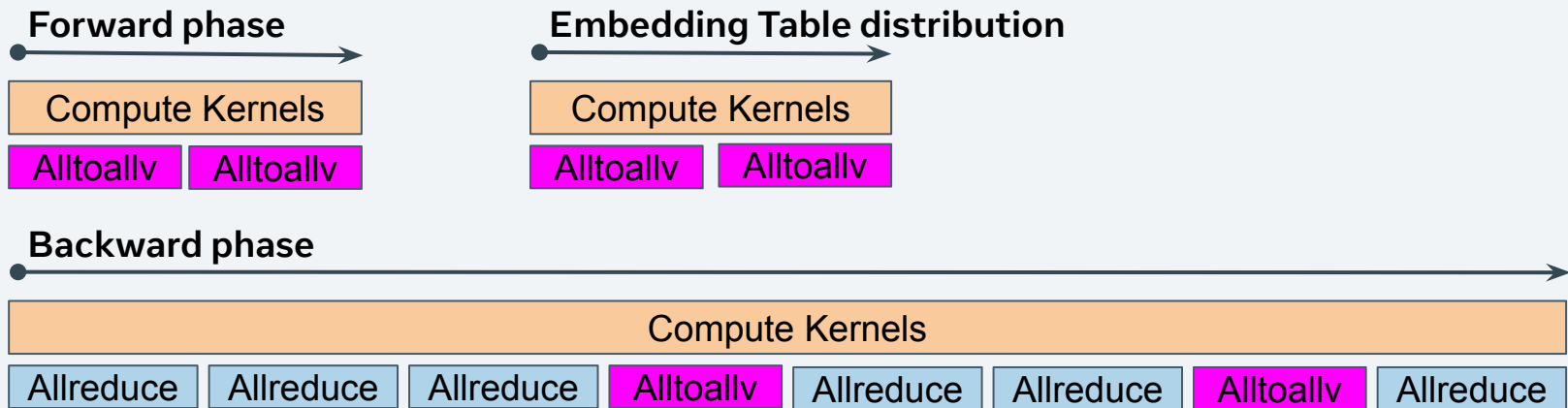
Distributed Training of DLRLMs

- Communication patterns depend on parallelism strategy

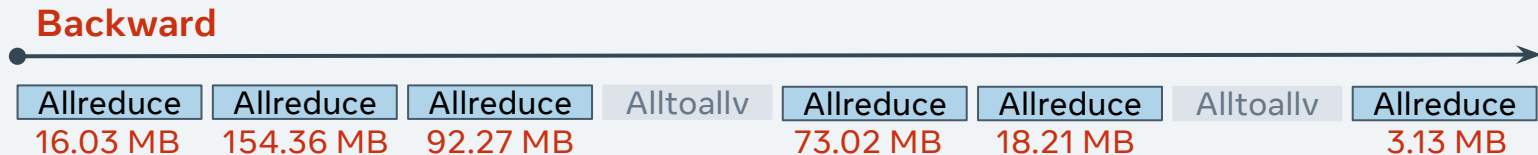


Communication in DLRM-2020 workloads

- **Key communication patterns**
 - **Allreduce** operations during backward phase → data parallelism
 - **Alltoallv** operations → model parallelism & table distribution
- **Message sizes and patterns are varied for different parallelisms**
 - Column-wise as an example



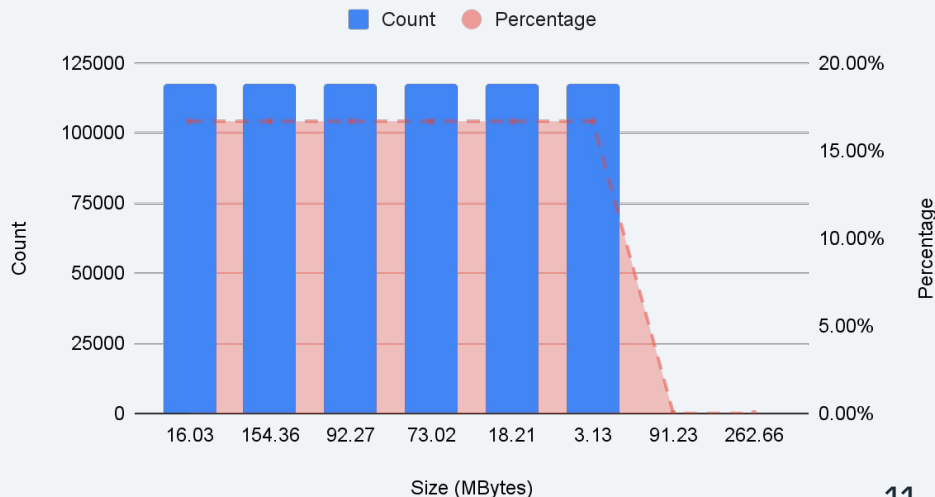
Allreduce Communication in DLRM-2020



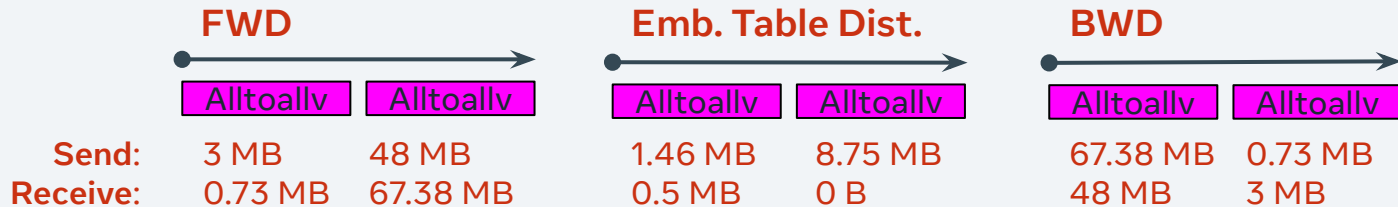
- **Variations***

- Batch size
- Model parallelism
 - Column-wise
- PyTorch parallelism
 - DDP
 - Bucket size
 - FSDP

Allreduce Count and Percentage in 128-GPU DLRM training



Alltoallv Communication in DLRM-2020

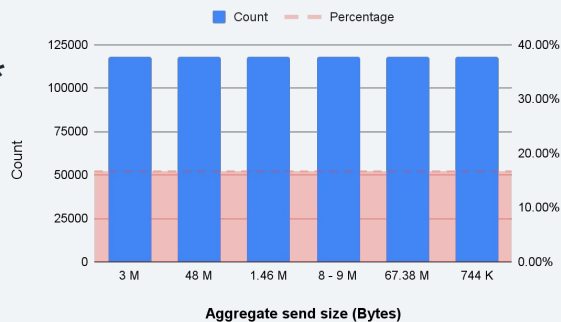


**These sizes are captured from rank-0, it is varied across ranks*

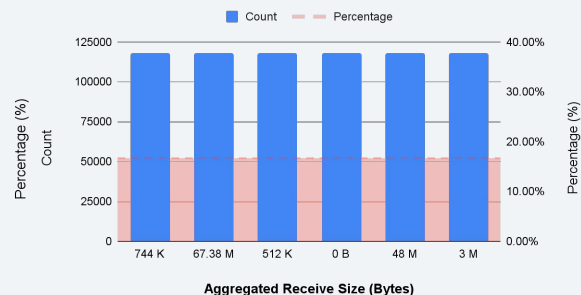
Variations

- Model Parallelisms
- Comms Quantization*

Alltoallv send sizes from Rank-0



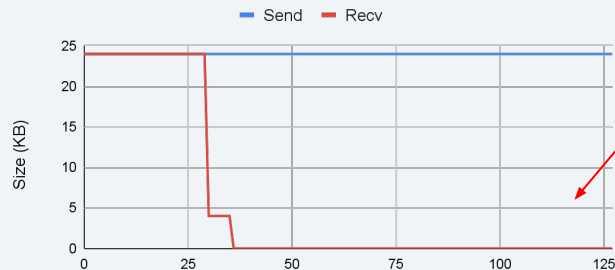
Alltoallv receive sizes from Rank-0



Alltoallv Communication in DLRM-2020 (cont.)

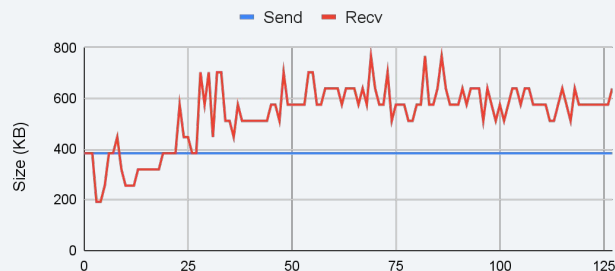
- Imbalanced communication

FWD Alltoallv-1



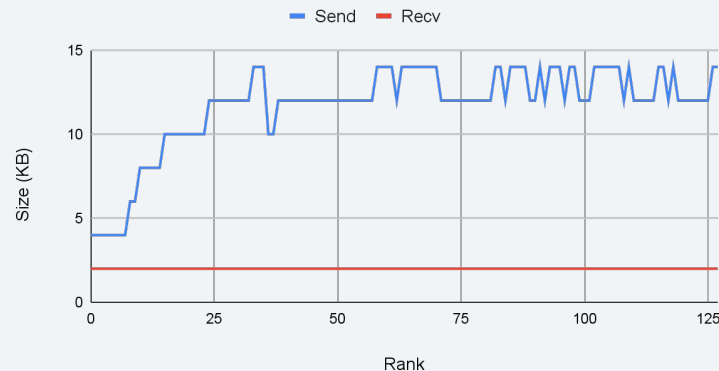
Zero data from most ranks

FWD Alltoallv-2

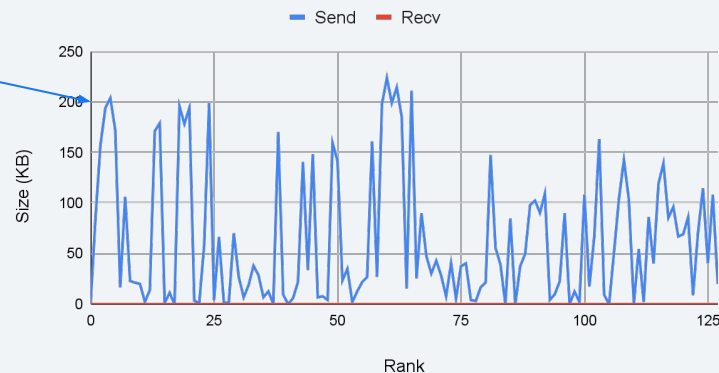


High skewness among ranks

Emb-Dist Alltoallv-1



Emb-Dist Alltoallv-2

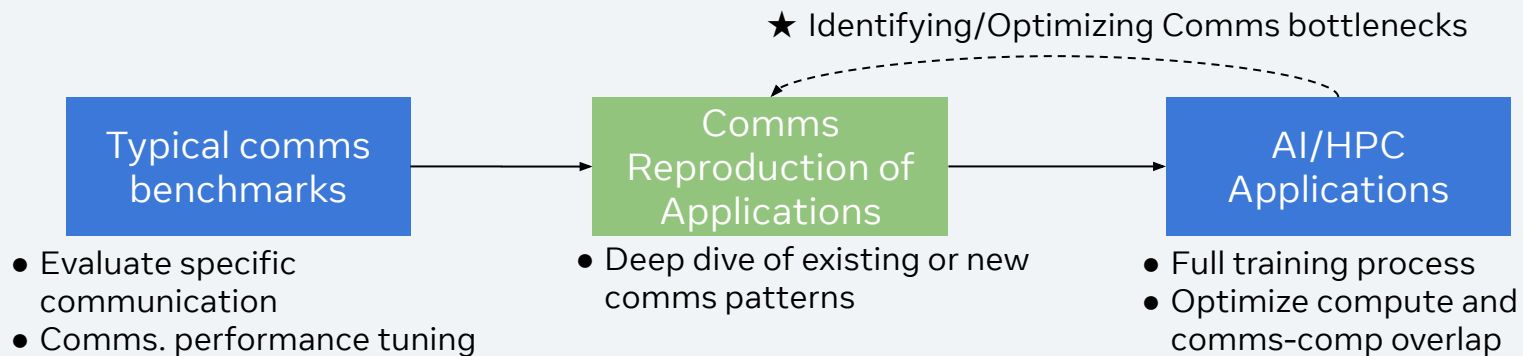


Agenda

- Introduction
- Demystify Communication Behavior in Training DLRM
- **Communication Reproduction**
 - Collect Communication trace
 - Replay communication trace
- Summary

Reproduction of Communication Behavior

- Why?
 - No need to run entire training workloads
 - Focus on understanding and optimizing communication patterns
 - No interference from computation
 - Cross-platform competitively analysis
 - E.g., Explore new SW/HW at scale for existing models



How to Replay Communication Trace

- Collecting communication traces
 - Captured from real-world production workloads
- Replaying communication traces using [PARAM benchmark](https://github.com/facebookresearch/param)
 - <https://github.com/facebookresearch/param>
 - PyTorch-based communication benchmark suite
 - Multi-backend support
 - NCCL, UCC, MPI, Gloo
 - Multi-device support
 - CPU, GPU (Nvidia & AMD), TPU

Summary

- Communication patterns in training DLRMs are complex, but predictable in general
 - Various parallelism methods, typically **~40%** time spent on communication
 - Large-sized Allreduce operations overlapped with compute kernels
 - Imbalanced Alltoallv operations are common
- Communication reproduction is important for SW/HW optimization
 - Open-source trace-replay benchmark (<https://github.com/facebookresearch/param>)
 - Need more communication traces of real-world workloads
- **Ongoing work**
 - Collecting, analyzing and sharing more communication traces of real-world DLRM workloads, e.g., various scales, various parallelism

