

### ZEN: Empowering Distributed Training with Sparsity-driven Data Synchronization

**Zhuang Wang**, Zhaozhuo Xu, Jingyi Xi, Yuke Wang Anshumali Shrivastava, T. S. Eugene Ng



## **Trends in Deep Learning**

- Data grows exponentially
  - Data Parallelism



Training dataset partitions

## **Trends in Deep Learning**

- Data grows exponentially
  - Data Parallelism
- Model size grows exponentially
  - Scale to trillion parameters
  - Model Parallelism: TP, PP, FSDP



### **Trends in Deep Learning Distributed training**

- Data grows exponentially
  - Data Parallelism lacksquare
- Model size grows exponentially
  - Scale to trillion parameters
  - Model Parallelism: TP, PP, FSDP lacksquare
- Hybrid parallelism for distributed training
  - Data Parallelism + Model Parallelism

### Gradient synchronization **Dense tensor synchronization among GPUs**

Communication and aggregation



- Synchronization of dense tensors has been well studied
  - BytePS<sup>[1]</sup> and Ring-AllReduce<sup>[2]</sup>

[1] A unified architecture for accelerating distributed DNN training in heterogeneous GPU/CPU clusters, OS [J] 2020 [2] Horovod: fast and easy distributed deep learning in TensorFlow, arXiv 2018

### **Sparsity in gradient tensors** Zero gradients in tensors

Gradients can be zeros

Low sparsity



High sparsity

### **Sparsity in gradient tensors** Zero gradients in tensors

- Gradients can be zeros
- Natural tensor sparsity in popular models

Model	MLP Size	Embedding Size	Embedding Sparsity
		88	8~ <b>F</b> ~~~~J
LSTM	20M	406M	98.87%
DeepFM	68M	214M	97.20%
NMT	31M	112M	97.53%

### **Sparsity in gradient tensors Zero gradients in tensors**

- Gradients can be zeros
- Natural tensor sparsity in popular models

Model	MLP Size	<b>Embedding Size</b>	<b>Embedding Sparsity</b>
LSTM	20M	406M	98.87%
DeepFM	68M	214M	97.20%
NMT	31M	112M	97.53%

Tensor sparsity from gradient sparsification algorithms <sup>[1,2,3]</sup>



Low sparsity

- [1] Sparse communication for distributed gradient descent, ACL 2017
- [2] Deep gradient compression: Reducing the communication bandwidth for distributed training, ICLR 2018
- [3] Sparsified SGD with Memory, NeurPS 2018

Up to 99% sparsity



### Synchronization of sparse tensors **Opportunities and Research Problem**

- Only communicate sparse tensors, i.e., non-zero gradients
  - Greatly reduces traffic volume and synchronization time













#### **Design space construction Four elemental dimensions**

Sparse tensor before synchronization



#### 2.5 GPU 3

# **Four elemental dimensions**

Sparse tensor before synchronization



Communication dimension



to-point	
hot	
lism	
ced	

# Four elemental dimensions

Sparse tensor before synchronization





Aggregation dimension



to-point	
hot	
	•
lism	
ced	

# **Four elemental dimensions**

Sparse tensor before synchronization



Aggregation dimension



to-point	
hot	
	•
lism	
ced	





to-point	
hot	
lism	
ced	

#### **Design space construction** Four elemental dimensions

Balance dimension



to-point
hot
lism
ced
*

#### **Design space construction Expressiveness**

Existing schemes can be described by the four dimensions

Schemes	Communication	Aggregation	Partition	Balance
AGsparse [1]	Ring, Hierarchy, Point-to-point	One-shot	Centralization	N/A
SparCML [2]	Hierarchy	Incremental	Centralization	N/A
OmniReduce [3]	Point-to-point	One-shot	Parallelism	Imbalanced

The four dimensions can describe the whole design space

- [1] Pytorch distributed: Experiences on accelerating data parallel, VLDB 2020
- [2] SparCML: High-performance sparse communication for machine learning, SC 2019
- [3] Efficient sparse collective communication and its application, SIGCOMM 2021

#### Find the optimal schemes **Based on the four dimensions**

• "Balanced Parallelism" (our proposal)





### Find the optimal schemes **Based on the four dimensions**

• "Balanced Parallelism" (our proposal)





[1] SparCML: High-performance sparse communication for machine learning, SC 2019



#### "Hierarchical Centraliza

ation"
to-point
hot
lism
ced

### Find the optimal schemes **Based on the four dimensions**

• "Balanced Parallelism" (our proposal)





[1] SparCML: High-performance sparse communication for machine learning, SC 2019



#### "Hierarchical Centraliza



ation"
to-point
hot
lism
ced

### How to realize "Balanced Parallelism"? **Imbalanced communication**





### How to realize "Balanced Parallelism"? **Challenge and requirements**

- Challenge: How to achieve balanced communications?
- Requirements

Support any workloads

> Balanced Parallelism

Preserve model accuracy

Small memory overhead

Small computation overhead

 Technique #1: Communication-oriented hash memory management Hash memory on GPU



 Technique #1: Communication-oriented hash memory management Hash memory on GPU

Parallel memory

9

8

Indices of non-zero gradients

2

5

12



 Technique #1: Communication-oriented hash memory management Hash memory on GPU



Technique #2: Multiple hash functions in each GPU thread

Hash memory on GPU



- Parallel hash

- ---- Rehash

Technique #2: Multiple hash functions in each GPU thread

Hash memory on GPU



Communication uncertainty due to rehashing

- Parallel hash
  - Serial write
- Communication
- Rehash

• Technique #3: Hierarchical consistent hashing across GPUs

Hash memory on GPU



First-level hash

Consistent across GPUs

- Parallel hash

- ---- Rehash

• Technique #3: Hierarchical consistent hashing across GPUs

Hash memory on GPU



First-level hash Second-level hash

Load-balanced communications with guarantee

- Parallel hash
  - Serial write
- Communication
- Rehash

### **Evaluations**

- Two testbeds
  - AWS EC2 p3dn.24xlarge with V100 GPUs, 100 Gbps RDMA networks
  - AWS EC2 p3.16xlarge with V100 GPUs, 25 Gbps TCP/IP networks
- Workloads
  - Models with natural tensor sparsity
  - Models with gradient compression (DGC <sup>[1]</sup>)

### **Communication** improvement 128 V100 GPUs

Communication speedups are normalized to AllReduce



Natural tensor sparsity

### **Communication** improvement 128 V100 GPUs

Communication speedups are normalized to AllReduce



### End-to-end training improvement 128 V100 GPUs

• Training throughput



![](_page_35_Figure_3.jpeg)

## Summary

- We study tensor sparsity in popular models
- We propose Zen with near-optimal synchronization for sparse tensors
  - find provably optimal schemes from the constructed design space
  - propose a hierarchical hashing algorithm for efficient realization on GPUs