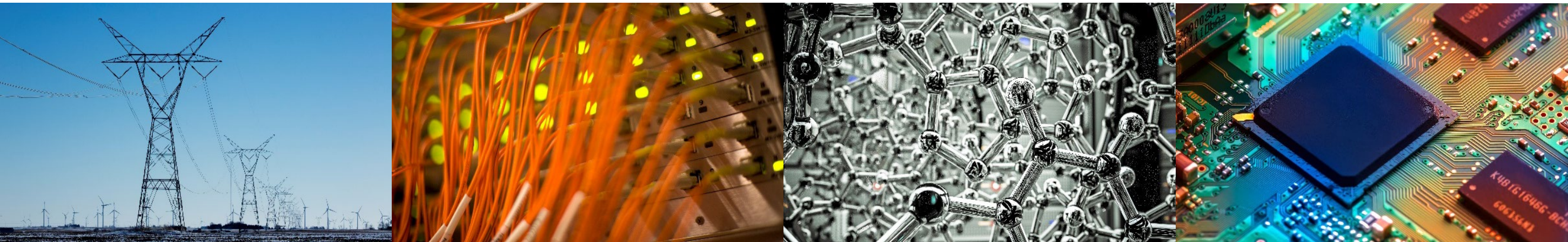


Large Graph Convolutional Network Training with GPU-Oriented Data Communication Architecture

Seung Won Min¹, Kun Wu¹, Sitao Huang¹, Mert Hidayetoğlu¹, Jinjun Xiong², Eiman Ebrahimi³, Deming Chen¹, Wen-mei Hwu^{1,3}

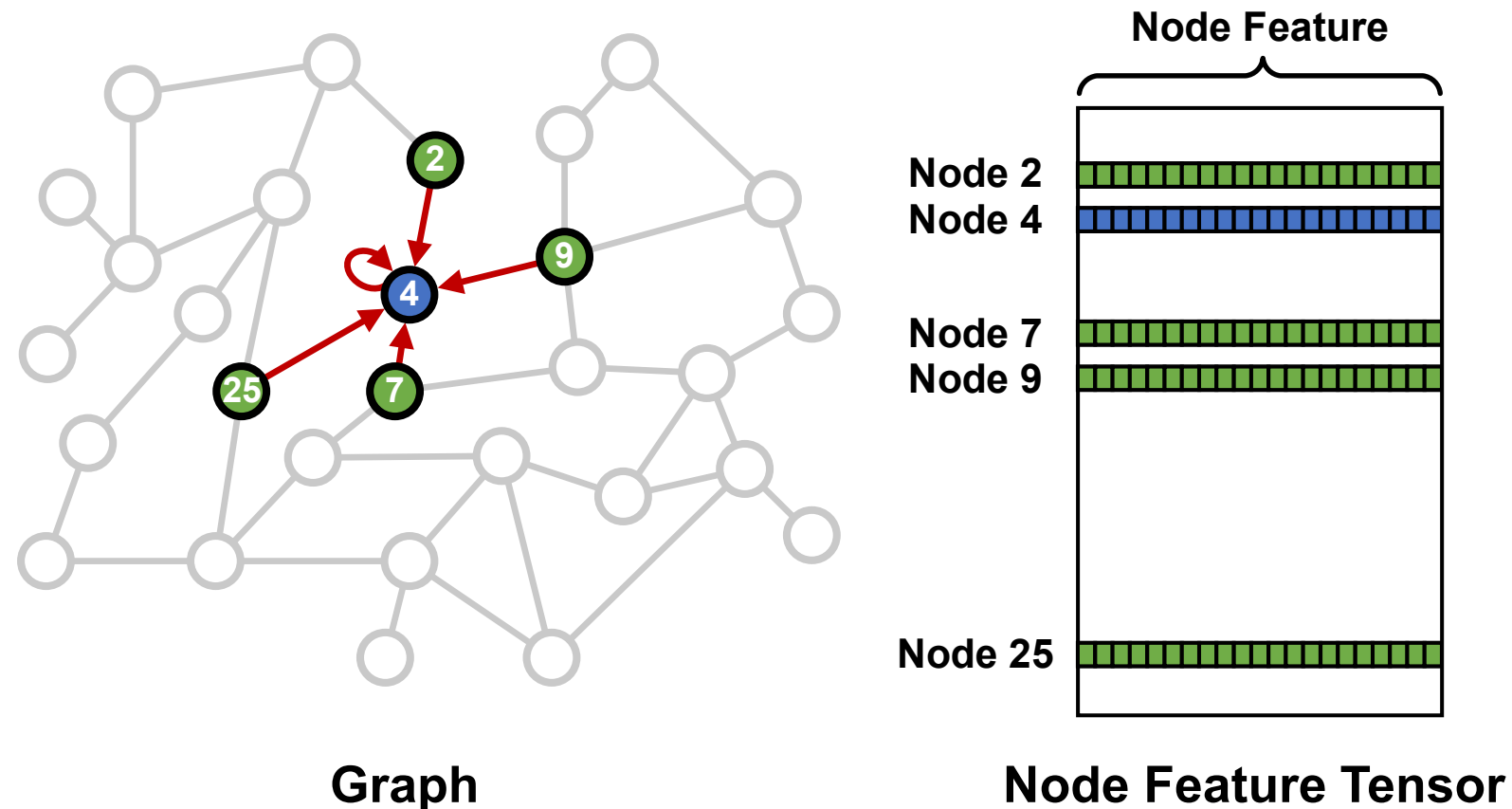
¹University of Illinois at Urbana-Champaign, ²IBM Research, ³NVIDIA Research

August 19th, 2021



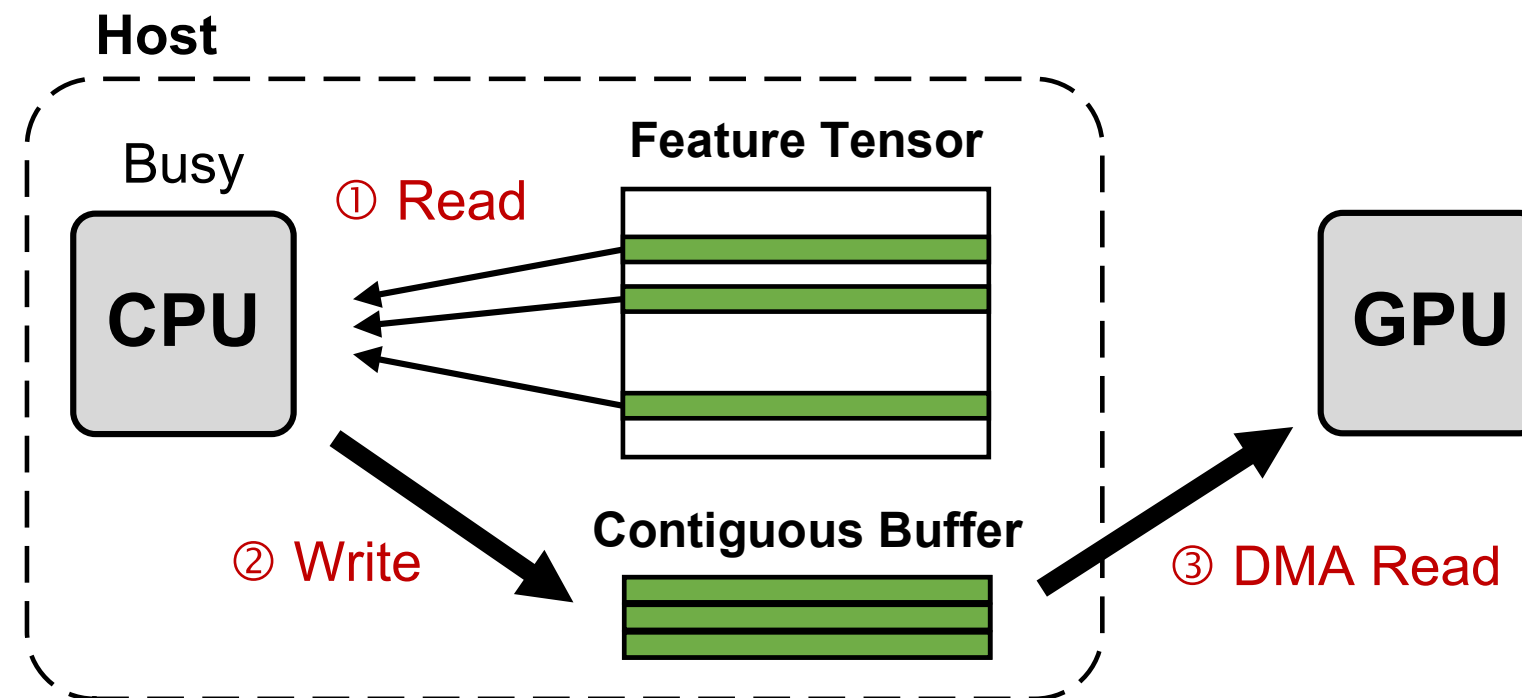
Graph Convolutional Network (GCN) Example

- Identify neighboring nodes and aggregate their features
- Problem statement: Transferring sparse data is not simple



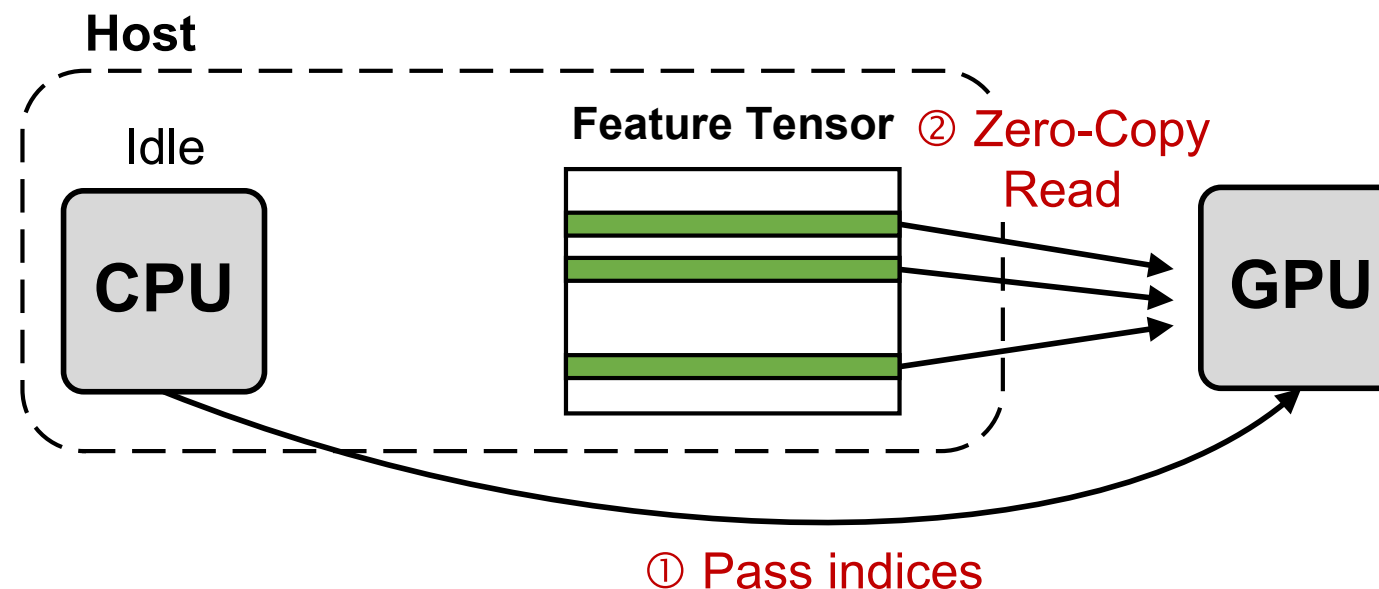
Current Method

- **Transform sparse data → dense data with CPU [1]**
 - Drawbacks: Long data access latency and waste of host resource
 - To sustain a single PCIe 4.0 x16 link, you consume $25 \times 3 = 75\text{GB/s}$



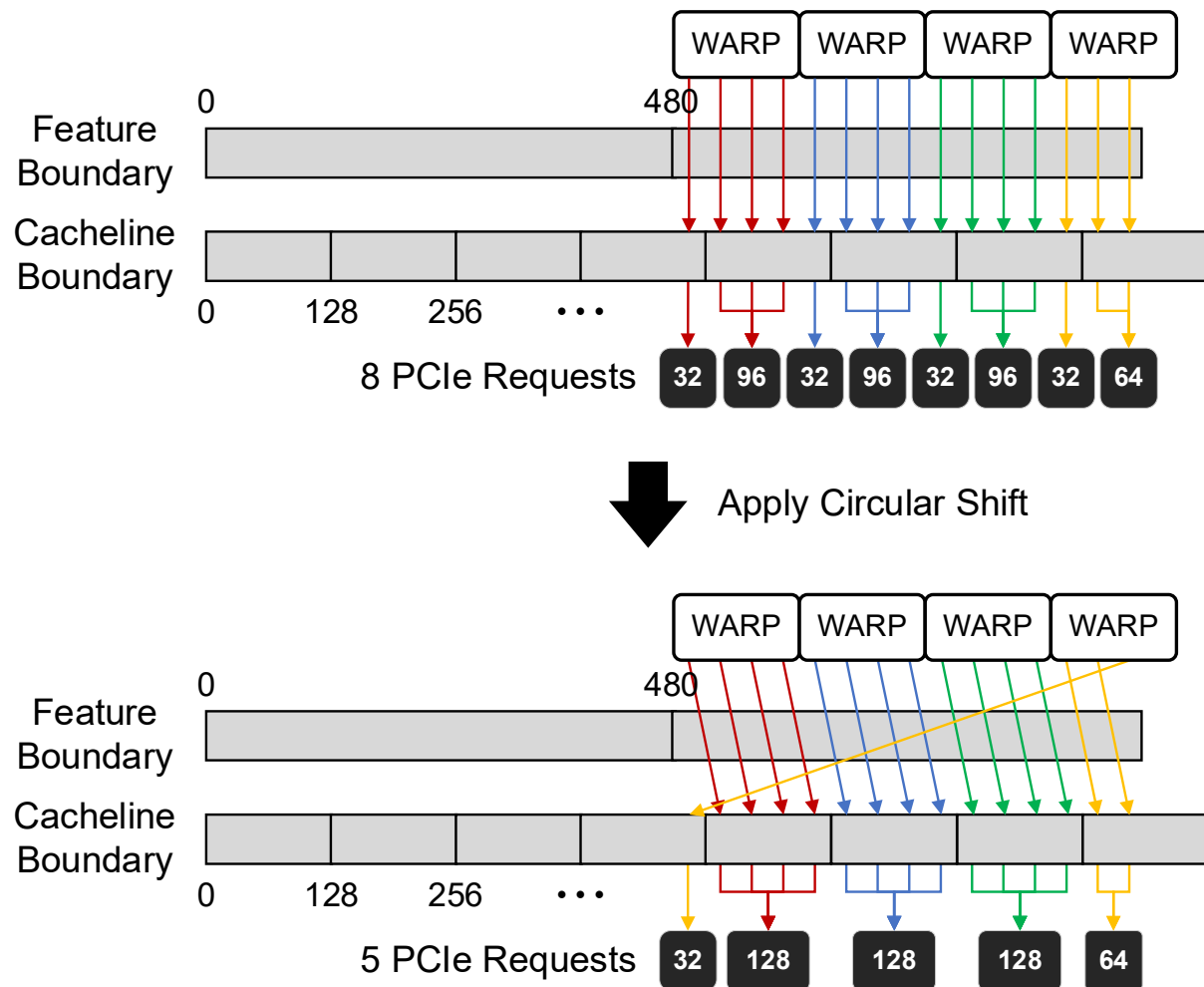
Proposed Method: Zero-Copy

- **GPUs directly access to the targeted node features**
 - Benefits: Less host resource usage (1/3 memory access) and reduced data access latency
 - Questions: Unclear actual performance benefits



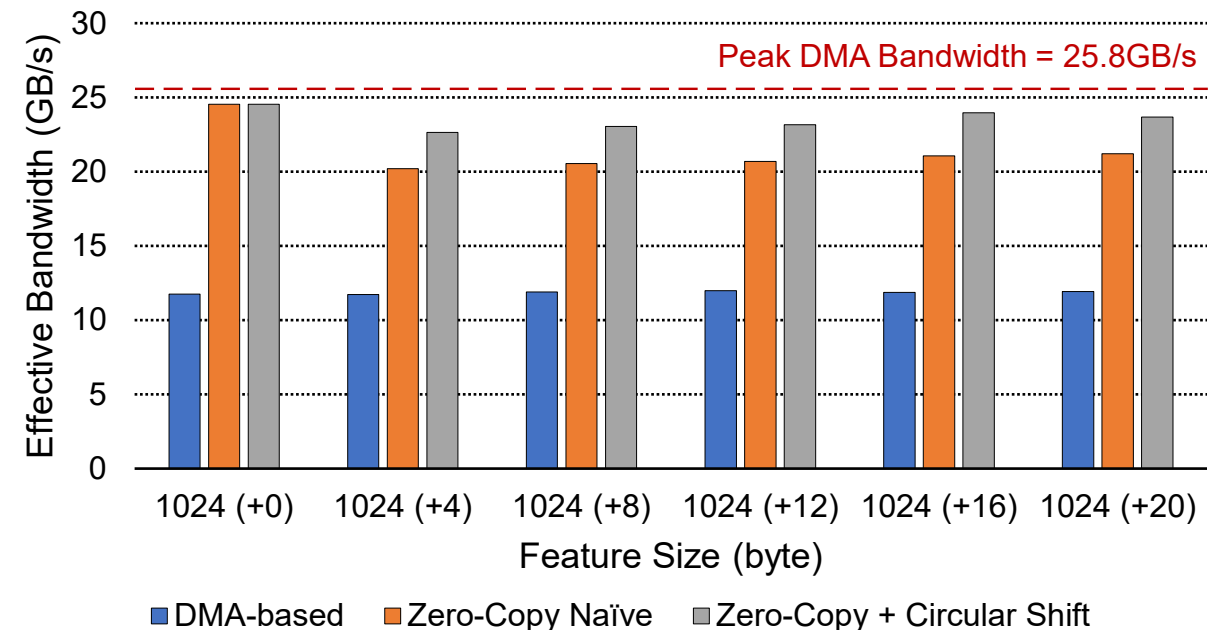
Zero-Copy Performance Optimization 1

Aligned PCIe memory read request



Automatic Alignment Adjustment

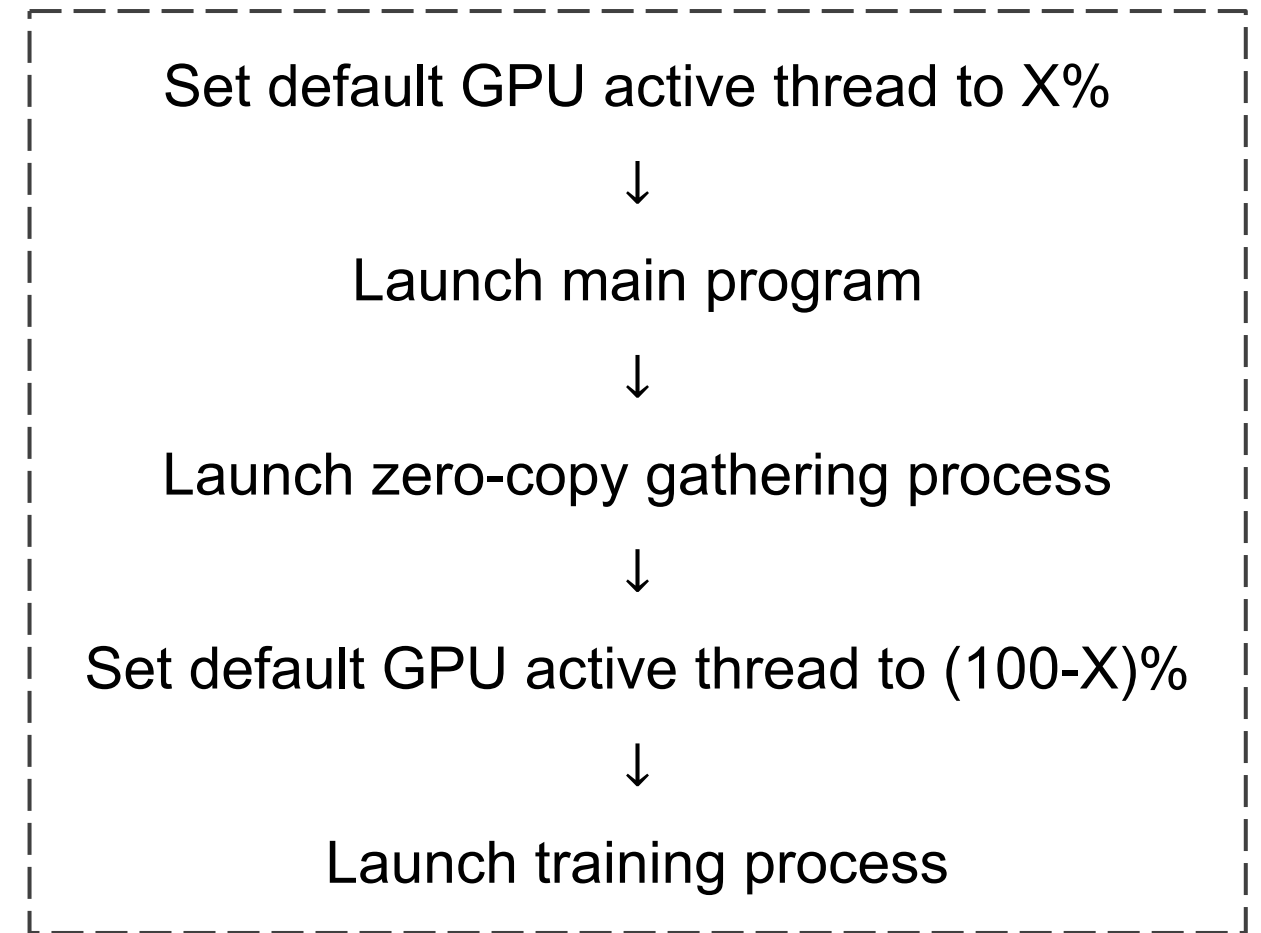
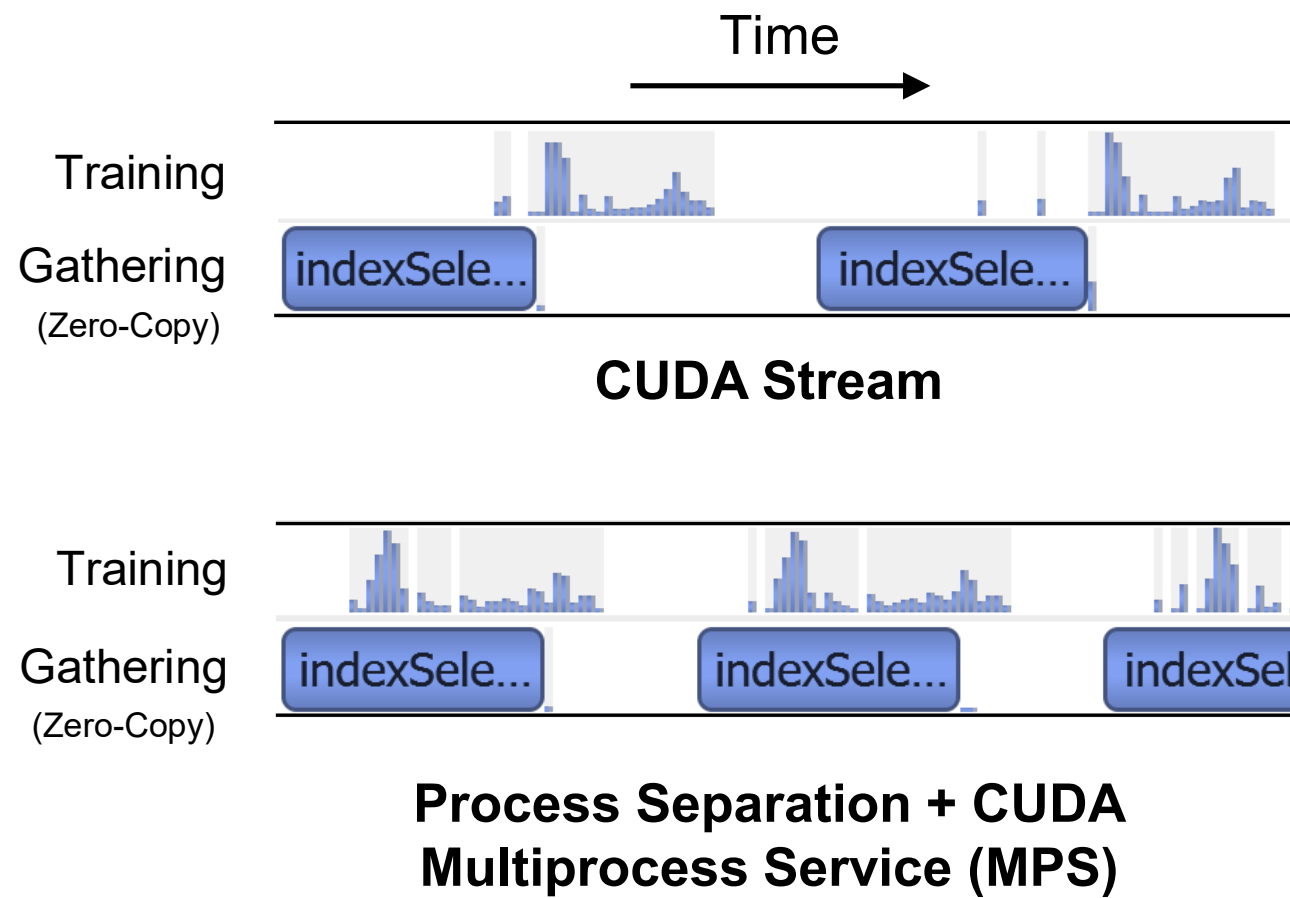
AMD Threadripper 3960x + NVIDIA RTX 3090



Sparse Feature Access Bandwidth

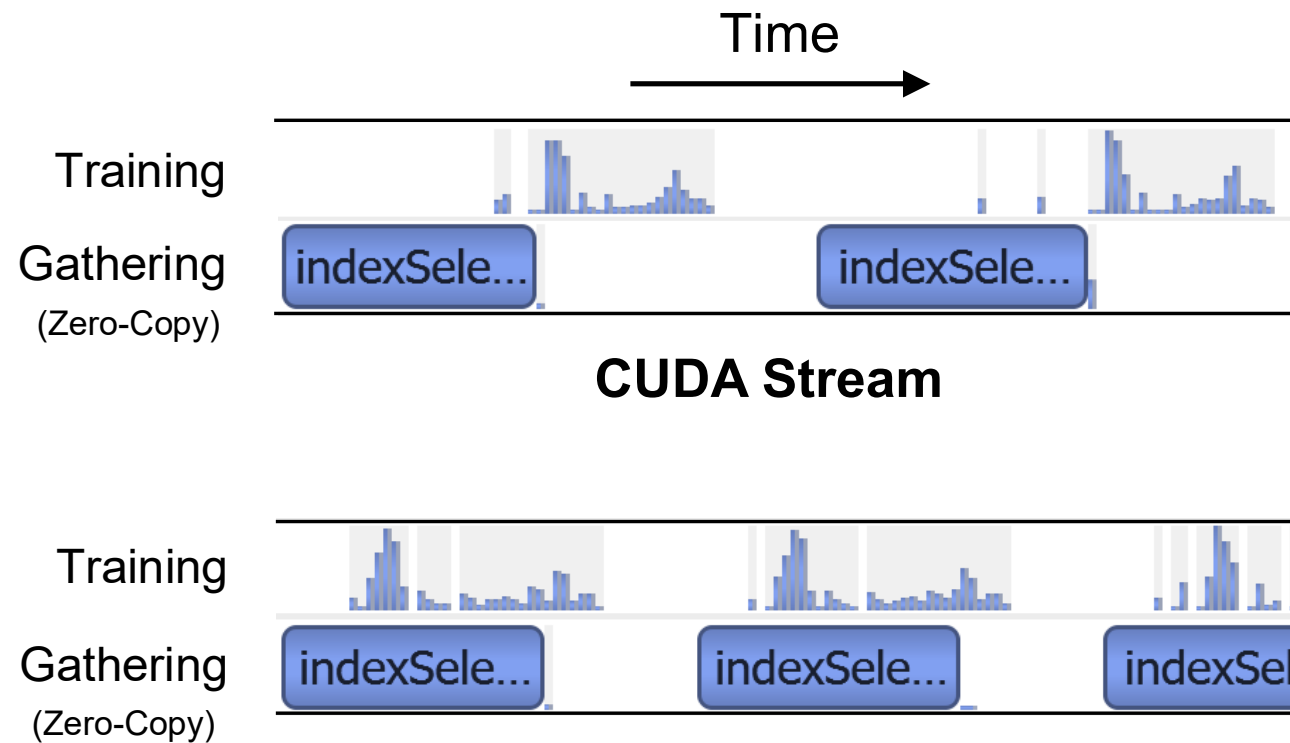
Zero-Copy Performance Optimization 2

Asynchronous zero-copy kernel operation



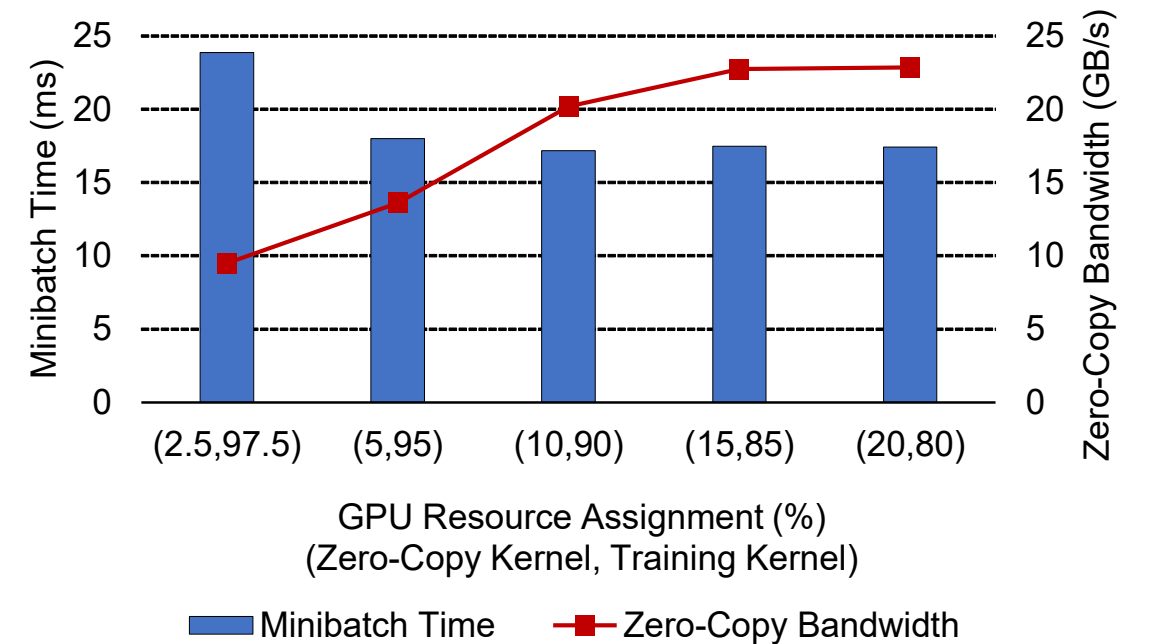
Zero-Copy Performance Optimization 2

Asynchronous zero-copy kernel operation



Process Separation + CUDA Multiprocess Service (MPS)

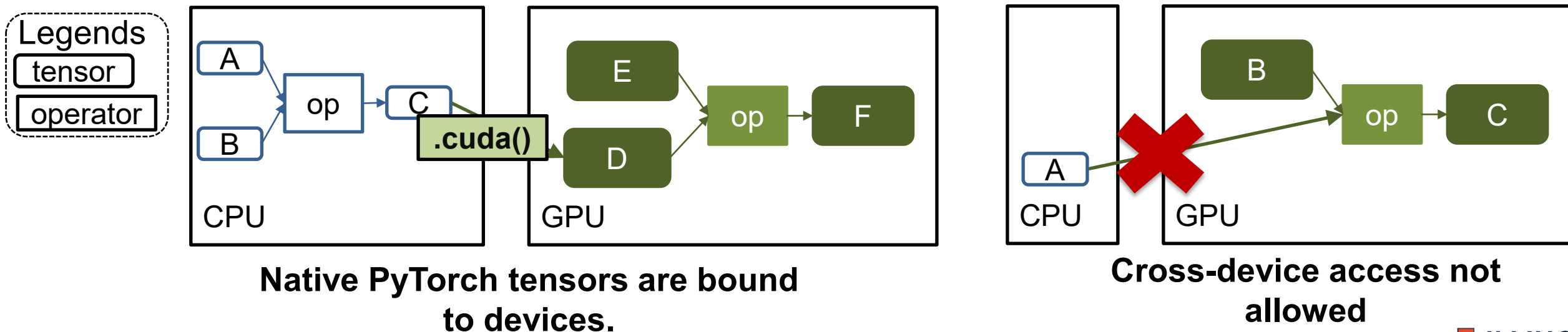
AMD Threadripper 3960x + NVIDIA RTX 3090



GraphSAGE Training Evaluation

Implementation: Major Challenges

- PyTorch 'always' keep data right next to the processing unit
 - Ex) CPU tensor in host memory and GPU tensor in GPU memory
 - No ways for GPU operators to directly access host memory
- Zero-copy tensor is transient
 - Result tensor should be GPU tensor in order to be stored in device memory.



Implementation: Unified Tensor

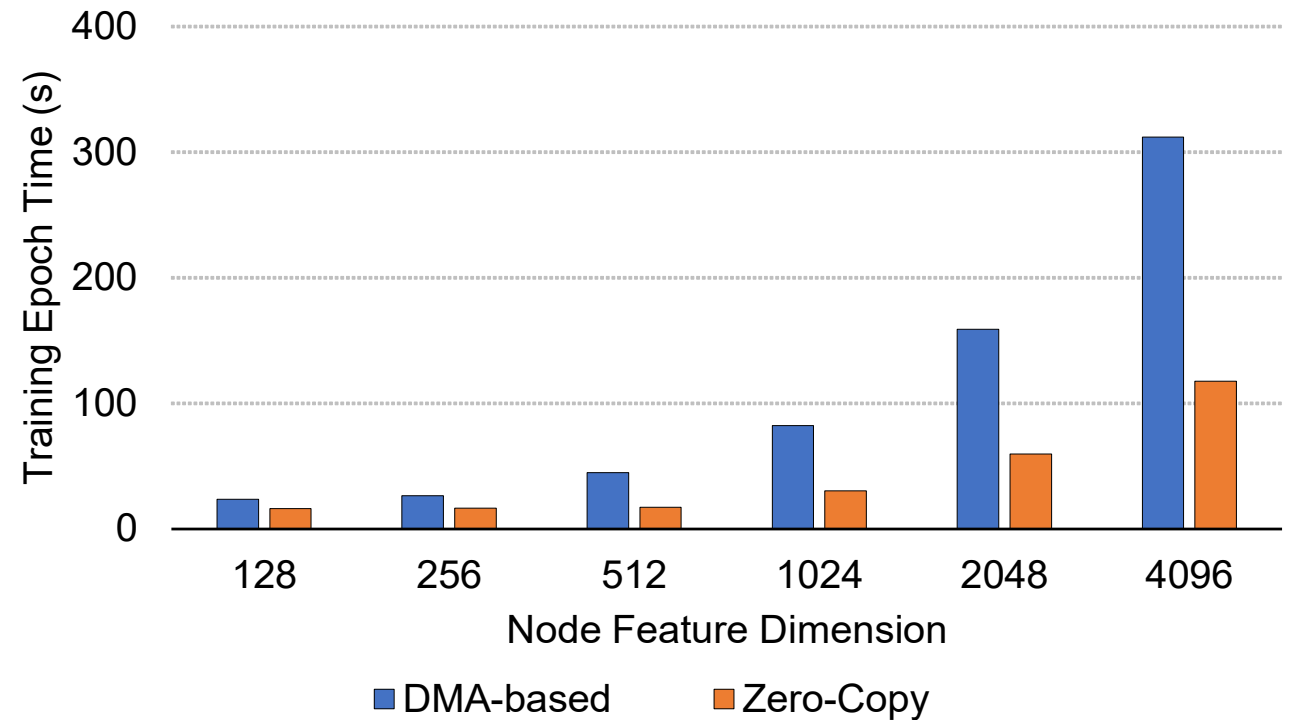
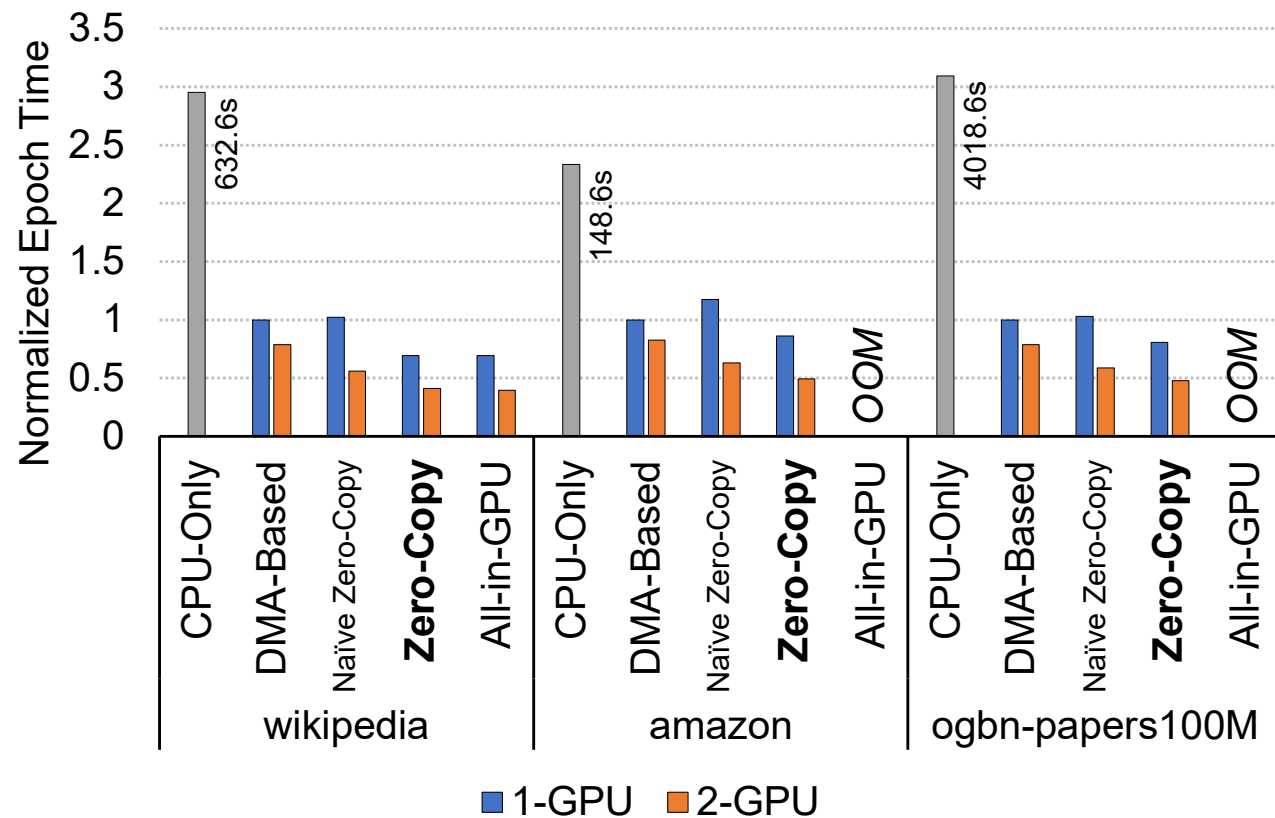
- We want a full-fledged new tensor type, to specify custom optimizations for zero-copy accesses
- Design of the “**Unified Tensor**”
 - Two affinity modes: GPU-affinitive and CPU-affinitive

Tensor type	Executor	Storage
CPU	CPU	CPU Memory
CUDA	GPU	GPU Memory
GPU-Affinitive Unified	GPU	GPU Memory
Host-Affinitive Unified	CPU	CPU Memory

- No actual data movements when the affinity is switched

Evaluations

- Deep Graph Library (DGL) used
- AMD Threadripper 3960x + NVIDIA RTX 3090 x2



Thank You!

Email: min16@illinois.edu

Github: https://github.com/K-Wu/pytorch-direct_dgl

DGL: <https://docs.dgl.ai/en/latest/api/python/dgl.contrib.UnifiedTensor.html>