# NeutronStar: **Distributed** GNN Training with **Hybrid** Dependency Management

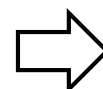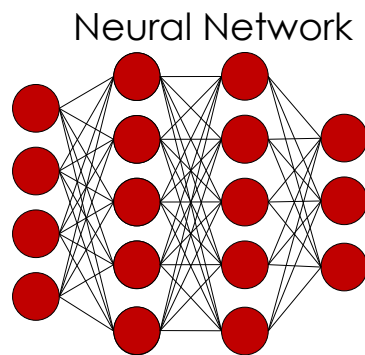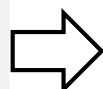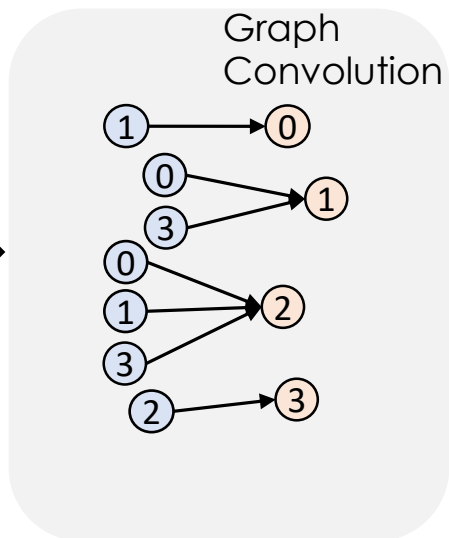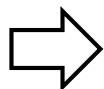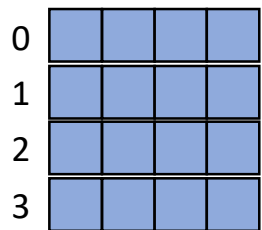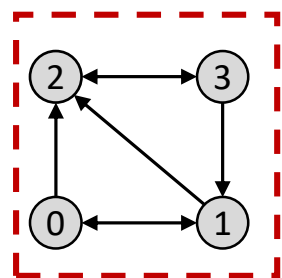**Qiange Wang**, Yanfeng Zhang, Hao Wang, Chaoyi Chen, Xiaodong Zhang,  Ge Yu

Northeastern University, China
International Digital Economy Academy (IDEA), China
The Ohio State University, USA
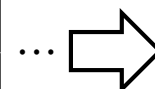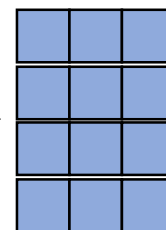
**SIGMOD 2022**

# Graph Neural Network



(a) Social Networks

(b) Knowledge Graph

(c) Biological networks

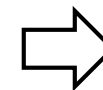# Graph Neural Network

Input data:



Embeddings

Predictions

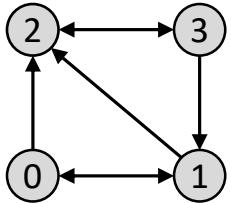Link Prediction

Node Classification

Protein Classification

# Execution Pattern of GNN

Forward computation (1-layer):

Input data:



| ScatterTo Edge | → | Edge Forward | → | Gather& Aggregate | → | Vertex Forward |
|---|---|---|---|---|---|---|

# Execution Pattern of GNN

Forward computation (1-layer):

Input data:

ScatterTo Edge → Edge Forward → Gather& Aggregate → Vertex Forward

Vertex 0:    i_Msg<0,2>:

Vertex 1:    i_Msg<1,2>:    Vertex 2:

Vertex 3:    i_Msg<3,2>:

Vertex Embedding[i]

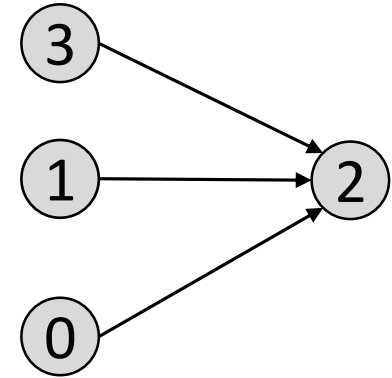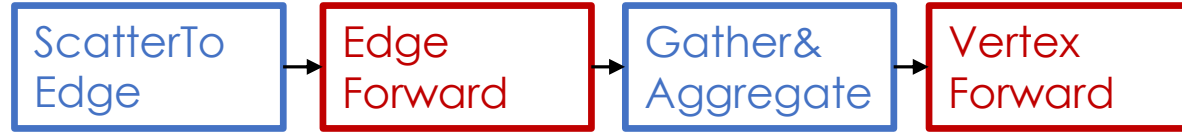# Execution Pattern of GNN

Forward computation (1-layer):

Input data:

ScatterTo Edge → Edge Forward → Gather& Aggregate → Vertex Forward



Vertex 0:    i_Msg<0,2>:    o_Msg<0,2>:

Vertex 1:    i_Msg<1,2>:    o_Msg<1,2>:

Vertex 3:    i_Msg<3,2>:    o_Msg<3,2>:

Vertex Embedding[i]
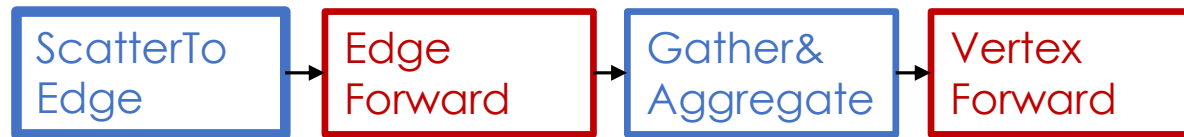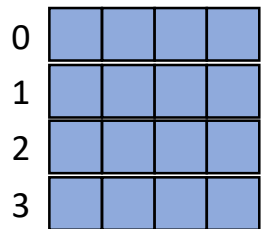
# Execution Pattern of GNN
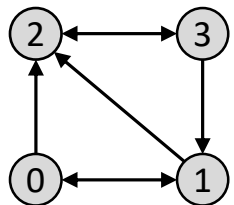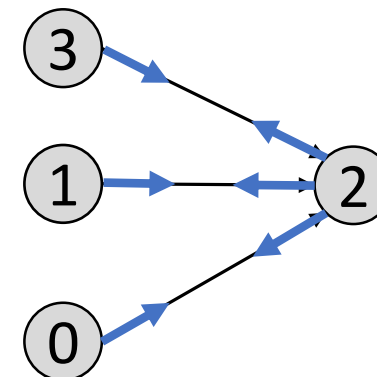
Forward computation (1-layer):

Input data:



ScatterTo Edge → Edge Forward → Gather& Aggregate → Vertex Forward

Vertex 0:   i_Msg<0,2>:   o_Msg<0,2>:

Vertex 1:   i_Msg<1,2>:   o_Msg<1,2>:   Vertex 2:

Neighbor Representation

Vertex 3:   i_Msg<3,2>:   o_Msg<3,2>:

Vertex Embedding[i]

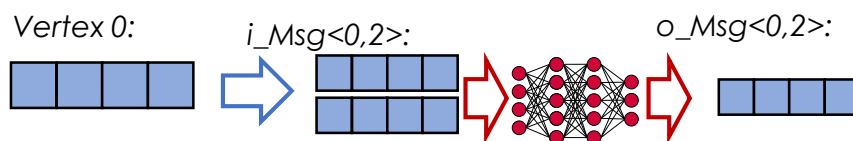# Execution Pattern of GNN
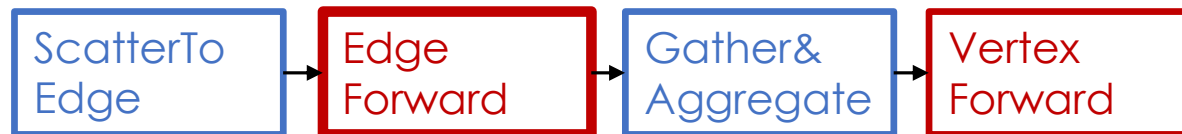
Forward computation (1-layer):

Input data:

ScatterTo Edge → Edge Forward → Gather& Aggregate → Vertex Forward

*Vertex 0:*  *i_Msg<0,2>:*  *o_Msg<0,2>:*

*Vertex 1:*  *i_Msg<1,2>:*  *o_Msg<1,2>:*

*Vertex 2:*  *Vertex 2:*

Neighbor Representation

Vertex Embedding[i+1]
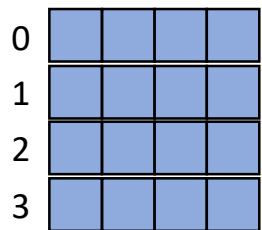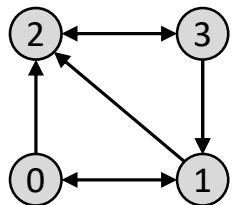
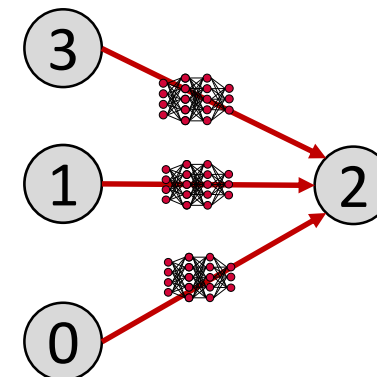*Vertex 3:*  *i_Msg<3,2>:*  *o_Msg<3,2>:*

Vertex Embedding[i]

# Execution Pattern of GNN

Forward computation (1-layer):

Input data:

ScatterTo Edge → Edge Forward → Gather& Aggregate → Vertex Forward
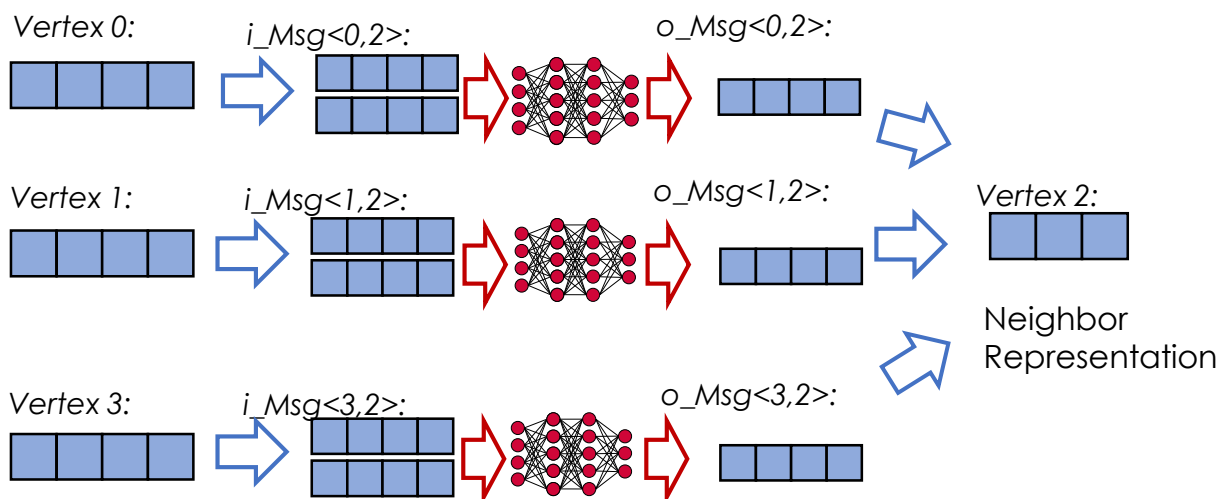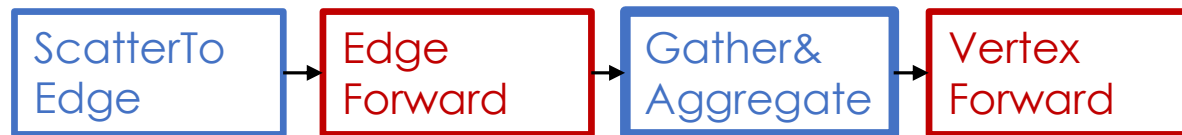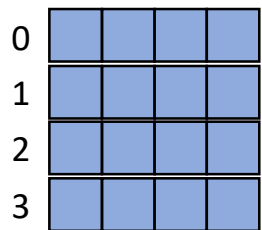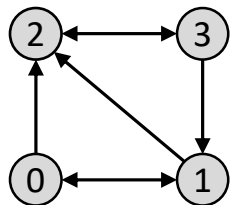


Vertex 0:

i_Msg<0,2>:

o_Msg<0,2>:

Vertex 1:

i_Msg<1,2>:

o_Msg<1,2>:

Vertex 2:

Neighbor Representation

Vertex 2:

Vertex Embedding[i+1]
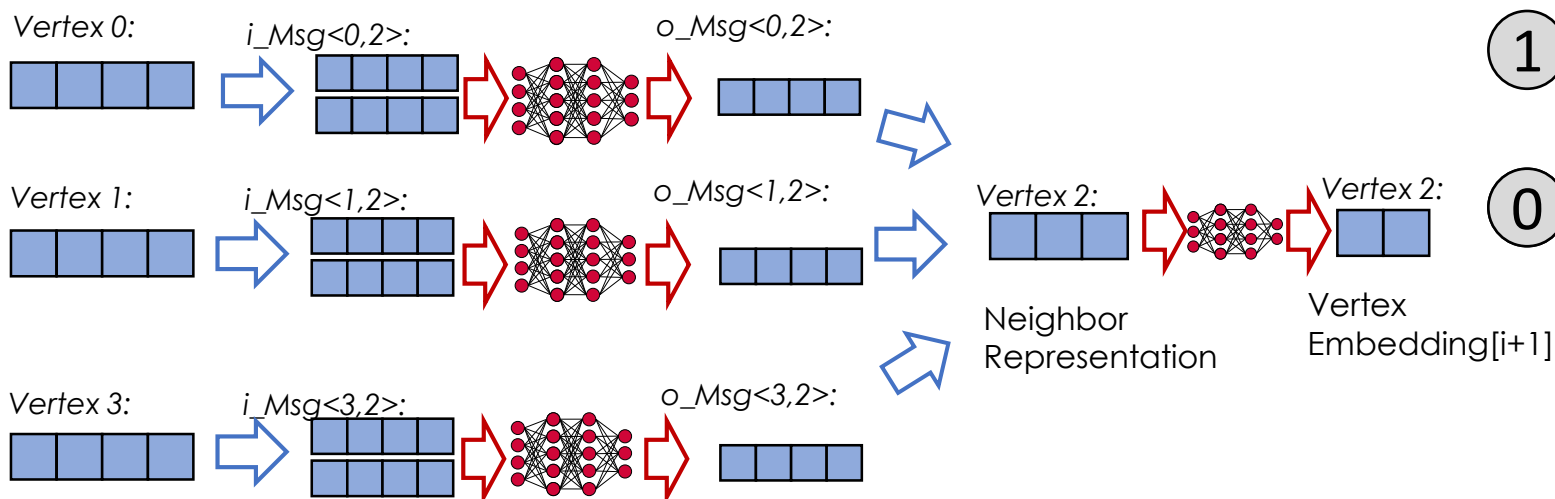
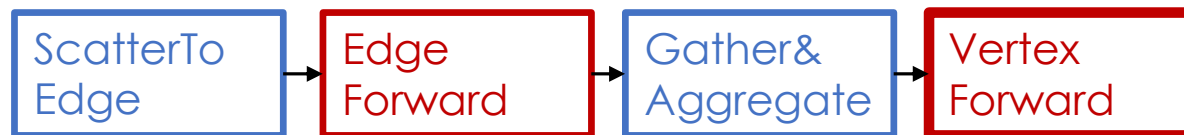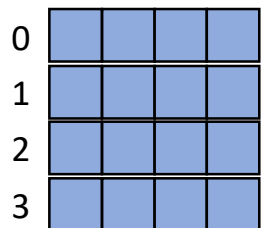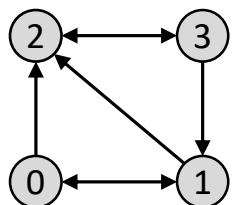**Downstream Computation**

Vertex 3:

i_Msg<3,2>:

o_Msg<3,2>:

Vertex Embedding[i]

# Execution Pattern of GNN

Forward computation:

Input data:

The computation of each vertex needs to gather information from its multi-hop neighbors

# Distributed GNN Training

DNN training:

Input data:



Proteins, Images, Sentences…

# Distributed GNN Training

DNN training:

Input data:

Proteins, Images, Sentences…

0
1
2
3

Partition 0

0
1

There is no cross-partition dependency

Partition 1

2
3

# Distributed GNN Training

DNN training:

Input data:

Proteins, Images, Sentences...

| 0 | | | | |
| 1 | | | | |
| 2 | | | | |
| 3 | | | | |

Partition 0

| 0 | | | | |
| 1 | | | | |

There is no cross-partition dependency

Partition 1

| 2 | | | | |
| 3 | | | | |

GNN training:

Input data:

| 0 | | | | |
| 1 | | | | |
| 2 | | | | |
| 3 | | | | |

# Distributed GNN Training

**DNN training:**

Input data:

Proteins, Images, Sentences…

Partition 0

There is no cross-partition dependency

Partition 1

**GNN training:**

Input data:

Node features, Graph topology

Graph Convolution

Graph Convolution

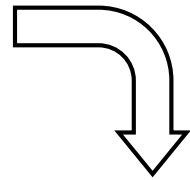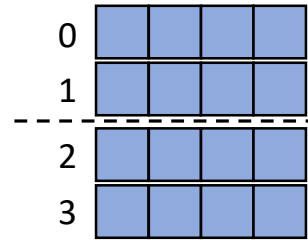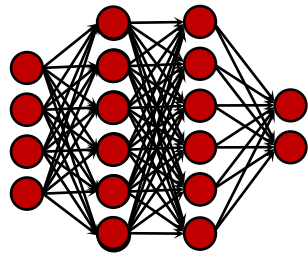Graph Convolution

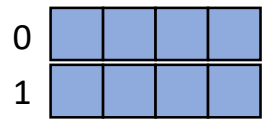Graph Convolution

# Distributed GNN Training
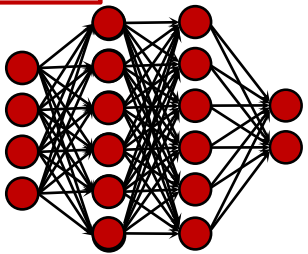
## DNN training:

Input data:

Proteins, Images, Sentences…
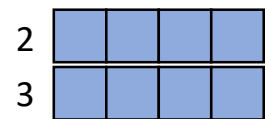
Partition 0

**There is no cross-partition dependency**

Partition 1

## GNN training:

Input data:

Node features, Graph topology

Graph Convolution

Graph Convolution

Graph Convolution

Graph Convolution

# Distributed GNN Training

Dependency tree of node 2:



2-hop neighbors      1-hop neighbors      Target vertex

GNN training:

Input data:

Node features, Graph topology

# Distributed GNN Training

Dependency tree of node 2:

GNN training:



Input data:

Node features,
Graph topology

Non-local dependencies
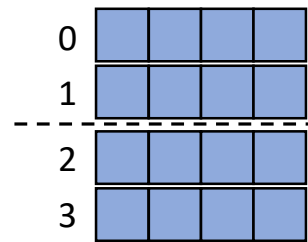
2-hop neighbors          1-hop neighbors          Target vertex

Graph
Convolution

Graph
Convolution

# Distributed GNN Training

Dependency tree of node 2:



Non-local dependencies

2-hop neighbors    1-hop neighbors    Target vertex

GNN training:

Input data:



Node features,
Graph topology

Distributed GNN tranining need to resolve the issues of vertex dependencies

# Challenges in Distributed Training

**Performance:**

Efficiently managing the cross-partition vertex representation.



2-hop neighbors      1-hop neighbors      Target vertex

# Challenges in Distributed Training

**Performance:**

Efficiently managing the cross-partition vertex representation.

**Usability:**

Automated cross-partition gradient backward propagation.



2-hop neighbors        1-hop neighbors        Target vertex

# Distributed GNN Training

Dependency tree of node 2:



2-hop neighbors      1-hop neighbors      Target vertex

Existing Approaches:

## **Dependency Cached:**

AliGraph[VLDB'20], Euler[arXiv'20], AGL[VLDB'20], DistDGL[arXiv'20]

# Dependency Cached Approach



Input data:

Dependency tree of Partition 0:

Dependency tree of Partition 1:

Remote Dependencies

2-hop neighbors    1-hop neighbors    Target vertices

Remote Dependencies

2-hop neighbors    1-hop neighbors    Target vertices

# Dependency Cached Approach

Input data:

Dependency tree of Partition 0:

Dependency tree of Partition 1:



2-hop neighbors    1-hop neighbors    Target vertices

2-hop neighbors    1-hop neighbors    Target vertices

# Dependency Cached Approach

Input data:

Cached dependencies of Partition 0:

Cached dependencies of Partition 1:

Remote Dependencies

Remote Dependencies

2-hop neighbors    1-hop neighbors    Target vertices

2-hop neighbors    1-hop neighbors    Target vertices

# Dependency Cached Approach

Input data:

Cached dependencies of Partition 0:

Cached dependencies of Partition 1:

**Remote Dependencies**

2-hop neighbors     1-hop neighbors     Target vertices

2-hop neighbors     1-hop neighbors     Target vertices

**Redundant computation problem**

# Distributed GNN Training

Dependency tree of node 2:



2-hop neighbors      1-hop neighbors      Target vertex

Existing Approaches:

**Dependency Communicated:**

ROC[MLSYS'20], Dorylus[OSDI'21], CAGNET[SC'20], DistGNN[SC'21], DGCL[EUROSYS'21].

# Dependency Communicated Approach



Input data:

Dependency tree of P0:

Dependency tree of P1:

Remote Dependencies

Remote Dependencies

2-hop neighbors  1-hop neighbors  Target vertices

2-hop neighbors  1-hop neighbors  Target vertices

# Dependency Communicated Approach

Input data:

Partitioned subgraph of P0:

Partitioned subgraph of P1:

# Dependency Communicated Approach

Input data:

Partitioned subgraph of P0:

Partitioned subgraph of P1:



Remote Dependencies

Remote Dependencies

2-hop neighbors    1-hop neighbors    Target vertices

2-hop neighbors    1-hop neighbors    Target vertices

# Dependency Communicated Approach

Input data:

Partitioned subgraph of P0:

Partitioned subgraph of P1:

**Remote Dependencies**

**Remote Dependencies**

2-hop neighbors    1-hop neighbors    Target vertices    2-hop neighbors    1-hop neighbors    Target vertices

# Dependency Communicated Approach

# Dependency Communicated Approach



Input data:

Partitioned subgraph of P0:

Partitioned subgraph of P1:

Remote Dependencies

Remote Dependencies

2-hop neighbors    1-hop neighbors    Target vertices

2-hop neighbors    1-hop neighbors    Target vertices

**Frequent cross-worker communication**

# Dependency Communicated Approach



Input data:

Partitioned subgraph of P0:

Partitioned subgraph of P1:

Remote Dependencies

Remote Dependencies

2-hop neighbors          1-hop neighbors          Target vertices          2-hop neighbors          1-hop neighbors          Target vertices

**Frequent cross-worker communication**

# Comparison of the Two Approaches

**Dependency Cached:**

**Dependency Communicated:**



The performance of **DepCache** and **DepComm** is dominated by the cost of (1) **redundant computation** and (2) **communication**.

# Comparison of the Two Approaches

**Dependency Cached:**



**Dependency Communicated:**



The performance of **DepCache** and **DepComm** is dominated by the cost of (1) **redundant computation** and (2) **communication**.

# Cost of the Two Approaches



2-hop neighbors      1-hop neighbors      Target vertex

## Cost of DepCache:

Graph convolution overhead:      Vertex computation overhead:



## Cost of DepComm:

Cross worker communication overhead

# Cost of the Two Approaches



2-hop neighbors          1-hop neighbors          Target vertex

(1) Graph inputs
(2) Model configurations
(3) Environment configurations

**Cost of DepCache:**

Graph convolution overhead:          Vertex computation overhead:



**Cost of DepComm:**

Cross worker communication overhead

# Comparison of the Two Approaches (1)

Graph inputs (the vertex degree):



(a) Graph Inputs

**Communication**

**Redundant Computation**

**DepComm is effective to high-degree vertices**

**DepCache is effective to low-degree vertices**

# Comparison of the Two Approaches (2)

Model configurations (Hidden layer sizes):



(b) Size of Hidden Layer

The volume of vertex dependencies increases exponentially with the number of hops

# Comparison of the Two Approaches (2)

Model configurations (Hidden layer sizes):



(b) Size of Hidden Layer

**Communication**

**Redundant Computation**

**DepComm is effective to large hidden layer size**

**DepCache is effective to small hidden layer size**

# Comparison of the Two Approaches (3)

Cluster Configurations (Computing power and network bandwidth):



(c) Cluster Environments

**DepComm is effective to high network bandwidth clusters**

**DepCache is effective to high computing power clusters**

# Hybrid Dependency Management



We evaluating the cost of **DepCache** and **DepComm** for each dependent neighbor.

# Hybrid Dependency Management

cost of **DepCache:**

$$t_r(v) = \sum_{l=0}^{L}\big(|E^l(v)| * T_e + Nbr^l(v)| * T_v\big) * dim^l$$

cost of **DepComm:**

$$t_c(v) = dim^l * T_c$$

Dependency tree of P1:



Remote Dependencies

2-hop neighbors     1-hop neighbors     Target vertices

We evaluating the cost of **DepCache** and **DepComm** for each dependent neighbor.

# Hybrid Dependency Management

cost of **DepCache:**



$$t_r(v) = \sum_{l=0}^{L} \left( |E^l(v)| * T_e + |Nbr^l(v)| * T_v \right) * dim^l$$

cost of **DepComm:**



$$t_c(v) = dim^l * T_c$$

Dependency tree of P1:



**Remote Dependencies**

2-hop neighbors      1-hop neighbors      Target vertices

We evaluating the cost of **DepCache** and **DepComm** for each dependent neighbor.

# Hybrid Dependency Management

COST MODEL:

$$T(P_i) = \sum_{v \in Dep_{Cache}} t_r(v) + \sum_{v \in Dep_{Comm}} t_c(v)$$

$$\text{s.t.,} \quad Size(Dep_{Cache}) < S$$

Dependency tree of P1:



**DepCache**

**Remote Dependencies**

**DepComm**

2-hop neighbors    1-hop neighbors    Target vertices

We evaluating the cost of **DepCache** and **DepComm** for each dependent neighbor.

# Hybrid Dependency Management

Input data:

Dependency tree of P0:

Dependency tree of P1:

**DepCache**

**Remote Dependencies**

**DepComm**

**Remote Dependencies**

**DepCache**

2-hop neighbors    1-hop neighbors    Target vertices

2-hop neighbors    1-hop neighbors    Target vertices

# Hybrid Dependency Management

Input data:

Dependency tree of P0:

Dependency tree of P1:



2-hop neighbors    1-hop neighbors    Target vertices      2-hop neighbors    1-hop neighbors    Target vertices

Remote Dependencies

# NeutronStar

We propose **NeutronStar**, a GPU-accelerated distributed GNN system with flexible automatic differentiation.

# Flexible Auto Differentiation

**Manually** implementing the **cross-worker** operators is **challenging**

# Flexible Auto Differentiation

**our goal:**

Implementing GNN computation with existing DNN libraries.

# Flexible Auto Differentiation

**our goal:**

Achieve efficient dependency management with graph engine.

# Flexible Auto Differentiation

Decoupling the **dependency management** and **GNN computation**



(a) Forward Computation

(b) Backward Computation

# Flexible Auto Differentiation

Decoupling the **dependency management** and **GNN computation**



Master-mirror communication

# Flexible Auto Differentiation

Decoupling the **dependency management** and **GNN computation**



Existing DL library (Pytorch) and NeutronStar's built-in graph operations.

# Flexible Auto Differentiation

Decoupling the **dependency management** and **GNN computation**



**Sync** - Compute

Master-to-Mirror
(Representation)

# Flexible Auto Differentiation

Decoupling the **dependency management** and **GNN computation**



Sync - **Compute**

Dependency Processing | Computation (Forward)

GetFromDep Neighbor | Forward One Layer

Computation (Backward) | Dependency Processing

Backward One Layer | PostToDep Neighbor

Forward Computation With Pytorch

# Flexible Auto Differentiation

Decoupling the **dependency management** and **GNN computation**



Backward Computation
with Pytorch

# Flexible Auto Differentiation

Decoupling the **dependency management** and **GNN computation**



Compute - **Sync**

Dependency Processing — Computation (Forward)

Computation (Backward) — Dependency Processing

Mirror-to-master (gradient of Representation)

# Experimental Setups

Baseline: **ROC**, **DistDGL**, **DepCache** (NeutronStar), and **DepComm** (NeutronStar).

Platforms:

    A 16-node Aliyun ECS cluster[1] (Each: 16 vCPUs, 62GB RAM, 1 NVIDIA-T4 GPU)

Algorithms and graphs:

- ☐ 3 Graph Neural Networks GCN, GIN, GAT
- ☐ 7 real-world graphs.

Environment

- ☐ Ubuntu 18.04 LTS
- ☐ CUDA 10.1

[1] Clusters are connected via 6GigE

**Table 2: Dataset description**

| Dataset | \|V\| | \|E\| | ftr. dim | #$\mathbb{L}$ | avg. deg | hid. dim |
|---|---|---|---|---|---|---|
| Google | 0.87M | 5.1M | 512 | 16 | 5.86 | 256 |
| Pokec | 1.6M | 30M | 512 | 16 | 18.75 | 256 |
| LiveJournal | 4.8M | 68M | 320 | 16 | 14.12 | 160 |
| Reddit | 0.23M | 114M | 602 | 41 | 487 | 256 |
| Orkut | 3.1M | 117M | 320 | 20 | 38.1 | 160 |
| Wiki-link | 12M | 378M | 256 | 16 | 31.12 | 128 |
| Twitter | 42M | 1.5B | 52 | 16 | 70.5 | 32 |

# Effectiveness of Hybrid Processing



Figure 11: Runtime results when varying the ratios between cached dependencies and communicated dependencies.

(1) Neither **communicating** nor **caching** all dependencies will reach the optimal performance.

(2) The optimal performance is reached when **mixing** **DepCache** and **DepComm**.

# Performance Comparison



**Figure 10: Overall performance comparison.**

Compared with the two representative distributed GNN systems **(DistDGL, ROC)**, **NeutronStar** achieves **1.8x – 14.3x** and **1.8X-5.3X** speedups on 3 GNNs and several real datasets, respectively.

# Performance Comparison



Figure 10: Overall performance comparison.

Compared with the **DepCache** and **DepComm**, **NeutronStar** achieves **2.0x – 15.0x** and **1.2X-1.7X** speedups on 3 GNNs and several real datasets, respectively.

# Accuracy Comparison

Time-to-accuracy comparison

NeutronStar outperforms other approaches

**1.20X** faster than DepComm
**1.96X** faster than DepCache-Sampling
**24.62X** faster than DepCache



Figure 14: Accuracy comparisons between Hybrid, DepComm, DepCache, **and DepCache-sampling with GCN on the Reddit dataset. Each dot indicates five training epochs for** Hybrid **and** DepComm, **and one training epoch for** DepCache **and DepCache-sampling.**

# Summary

NeutronStar: Distributed GNN training with hybrid dependency management.

☐ **Providing insight into the two existing approaches**

We conduct a comprehensive study on the performance merits and limits of the two distributed GNN training approaches (**DepCache** and **DepComm**).

# Summary

NeutronStar: Distributed GNN training with hybrid dependency management.

☐ **Providing insight into the two existing approaches**

We conduct a comprehensive study on the performance merits and limits of the two distributed GNN training approaches (**DepCache** and **DepComm**).

☐ **Proposing a hybrid dependency management framework**

We identify the key tradeoff between the two approaches and propose a hybrid dependency management approach.

# Summary

NeutronStar: Distributed GNN training with hybrid dependency management.

☐ **Providing insight into the two existing approaches**

We conduct a comprehensive study on the performance merits and limits of the two distributed GNN training approaches (**DepCache** and **DepComm**).

☐ **Proposing a hybrid dependency management framework**

We identify the key tradeoff between the two approaches and propose a hybrid dependency management approach.

☐ **Delivering a fast distributed GNN system**

We design and implement NeutronStar, a distributed GNN system with auto differentiation that achieves 1.8X-14.3X speedups over the existing GNN systems.

# Summary

NeutronStar: Distributed GNN training with hybrid dependency management.

- **Providing insight into the two existing approaches**
  We conduct a comprehensive study on the performance merits and limits of the two distributed GNN training approaches (**DepCache** and **DepComm**).

- **Proposing a hybrid dependency management framework**
  We identify the key tradeoff between the two approaches and propose a hybrid dependency management approach.

- **Delivering a fast distributed GNN system**
  We design and implement NeutronStar, a distributed GNN system with auto differentiation that achieves 1.8X-14.3X speedups over the existing GNN systems.

- **The codes are publicly available on github**

  https://github.com/Wangqge/NeutronStarLite

Questions