

2025 USENIX Annual Technical Conference

Accelerating distributed graph learning by using collaborative in-network multicast and aggregation

Zhaoyi Li^{1,3}, Jiawei Huang¹, Yijun Li¹, Jingling Liu¹, Junxue Zhang², Hui Li¹
Xiaojun Zhu¹, Shengwen Zhou¹, Jing Shao¹, Xiaojuan Lu¹, Qichen Su¹
Jianxin Wang¹, **Chee Wei Tan**³, Yong Cui⁴, Kai Chen²

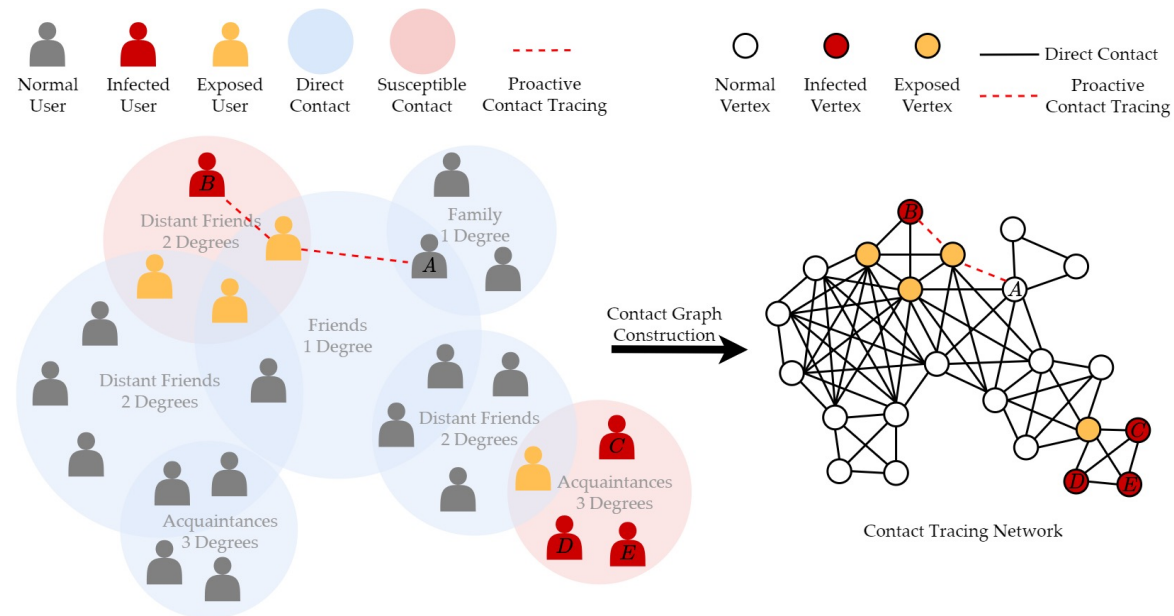
¹Central South University, ²Hong Kong University of Science and Technology

³Nanyang Technological University, ⁴Tsinghua University

Introduction

■ Large-scale Computation of Graph Neural Networks

- Drug discovery, public health surveillance, online social networks

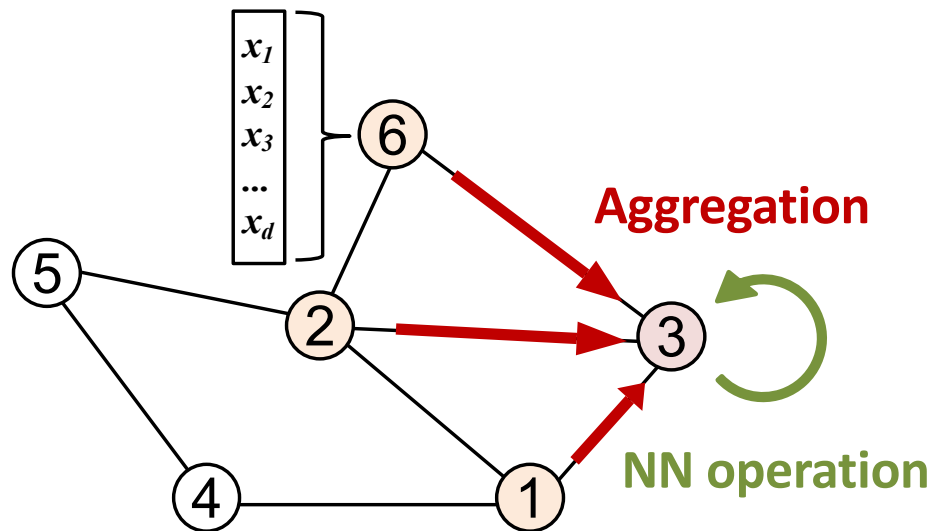


Chee Wei Tan et al, *DeepTrace: Learning to optimize contact tracing in epidemic networks with graph neural networks*, IEEE Trans. on Signal and Information Processing over Networks, 11:97–113, 2025.

Introduction

■ Graph Neural Network

- GNN can capture structured information from graph-based data



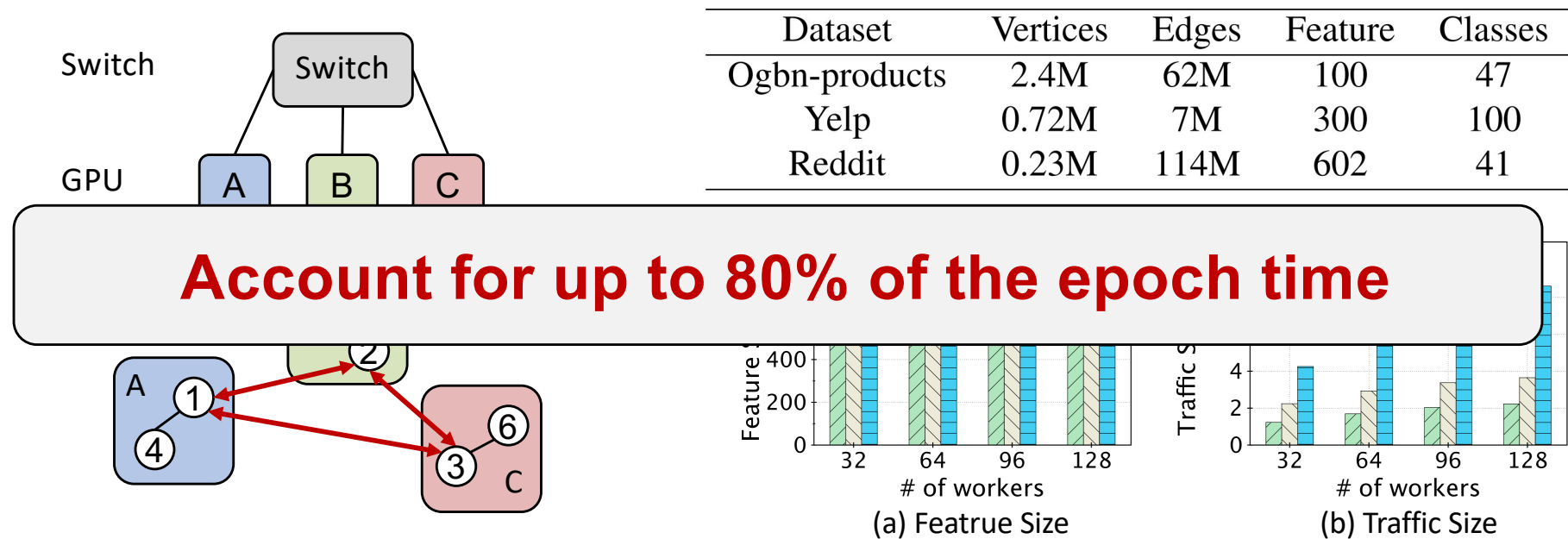
Basic Steps:

- (1) Neighbor aggregation
- (2) Vertex representation update

Introduction

■ Communication Bottleneck in Full-Graph Training

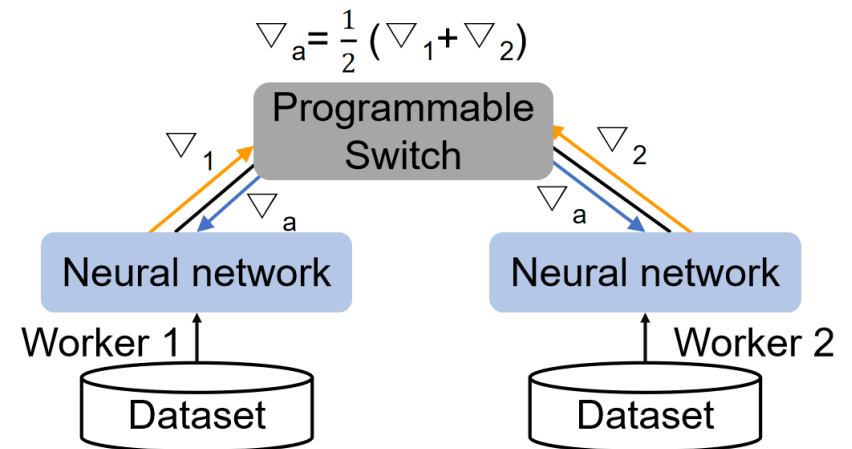
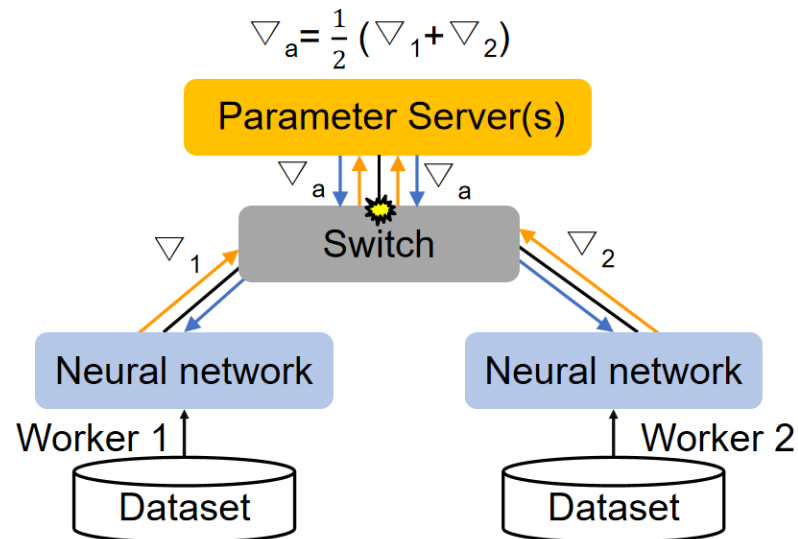
- Host-based multicast and aggregation leads to 1-to-N redundant traffic and N-to-1 bandwidth contention



Introduction

■ In-Network Aggregation (INA)

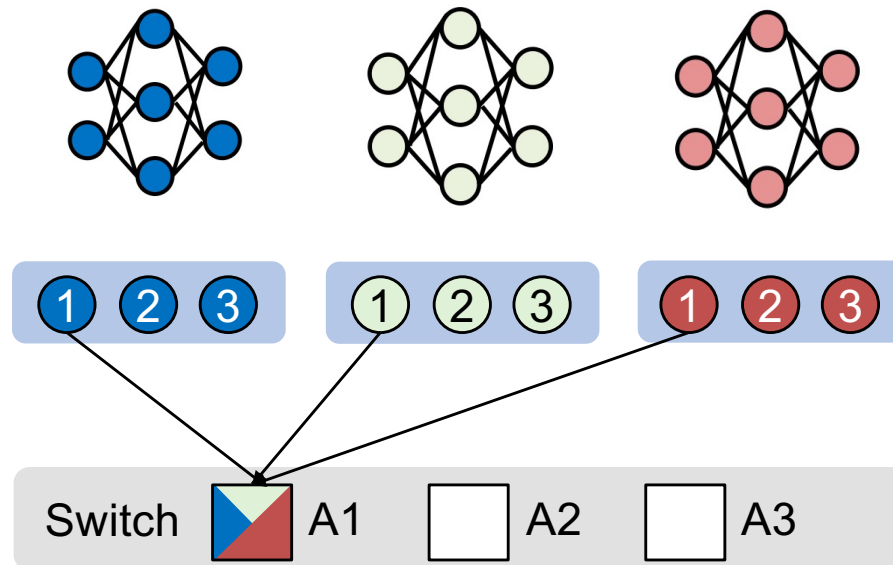
- Offloading the gradient aggregation into programmable switch
- Reducing traffic volume for data-parallel training



Introduction

■ In-Network Aggregation SwitchML [NSDI'21], ATP [NSDI'21]

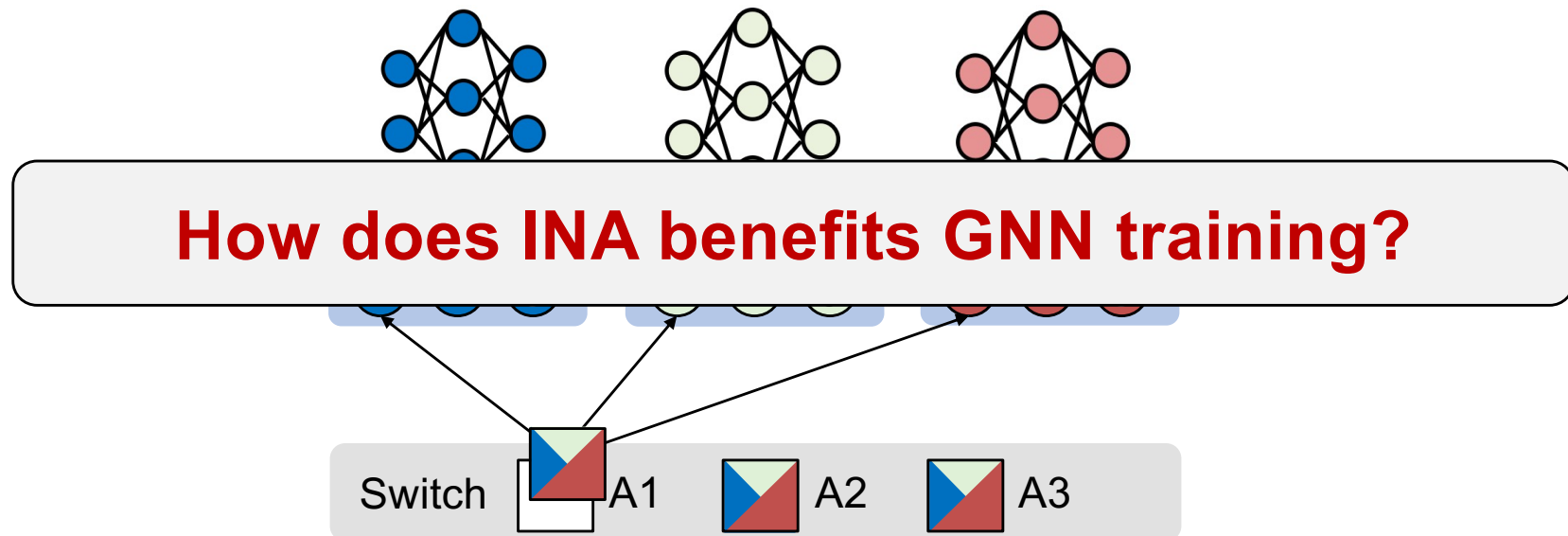
- Offloading the gradient aggregation into programmable switch
- Reducing traffic volume for data-parallel training



Introduction

■ In-Network Aggregation SwitchML [NSDI'21], ATP [NSDI'21]

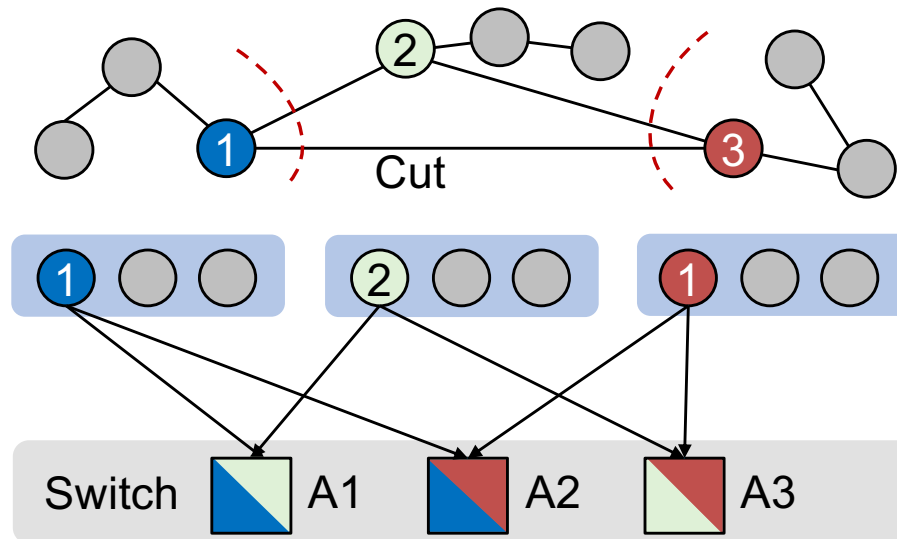
- Offloading the gradient aggregation into programmable switch
- Reducing traffic volume for data-parallel training



Motivation

■ Strawman Solution

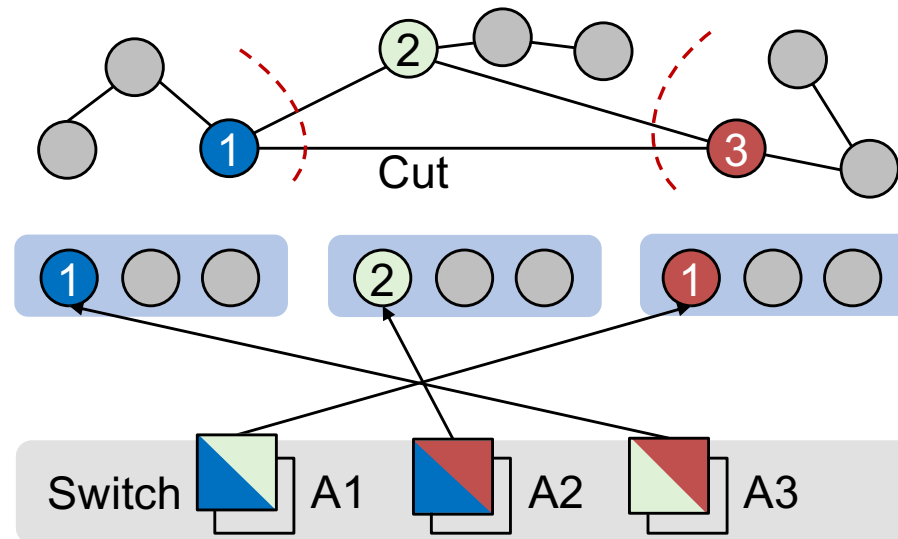
- Offloading multicast and aggregation of vertices feature into switches
- Eliminating 1-to-N redundant traffic and N-to-1 bandwidth bottleneck



Motivation

■ Strawman Solution

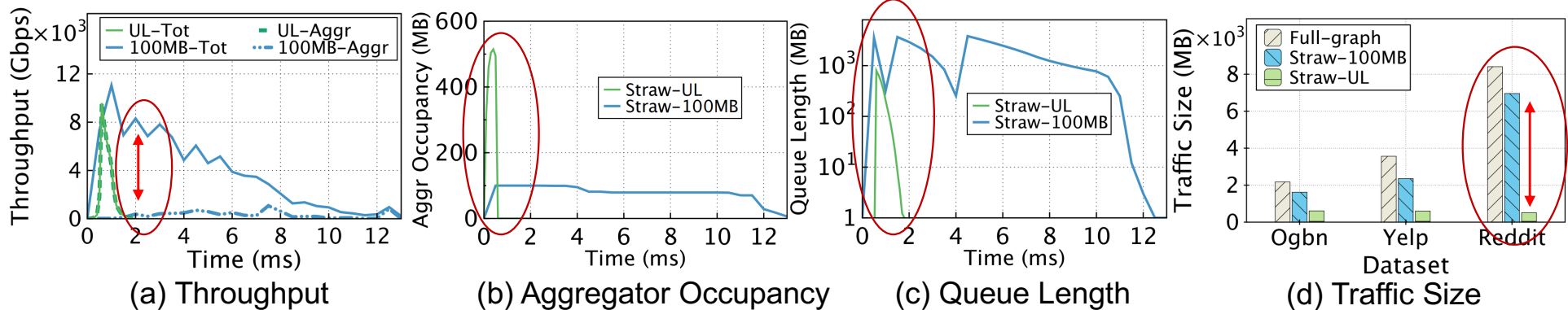
- Offloading multicast and aggregation of vertices feature into switches
- Eliminating 1-to-N redundant traffic and N-to-1 bandwidth bottleneck



Motivation

■ Experimental Observations

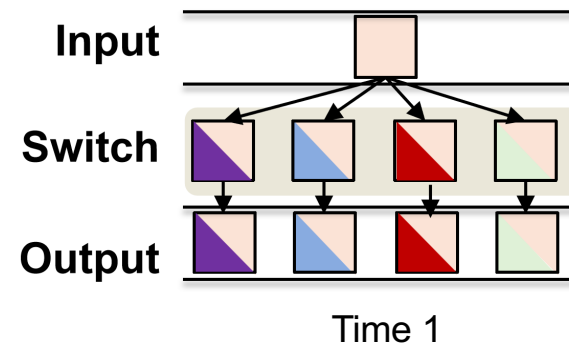
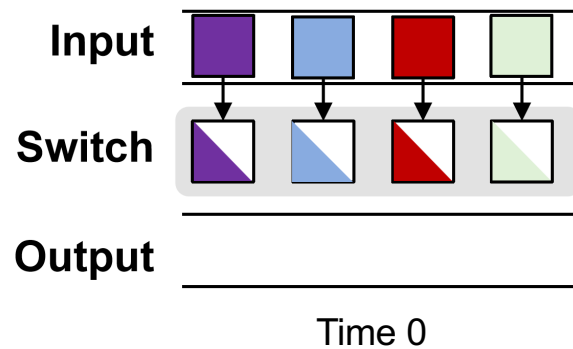
- Evaluating strawman solution under limited switch memory (100MB)
- Suffering from low aggregation throughput, significant queue backlog, and large traffic volume



Motivation

■ Deep Dive

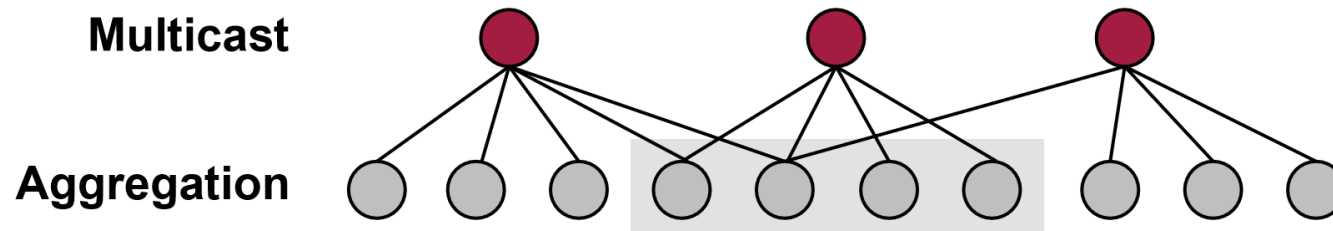
- Graph-agnostic multicast order results in inefficient forwarding pipeline



Motivation

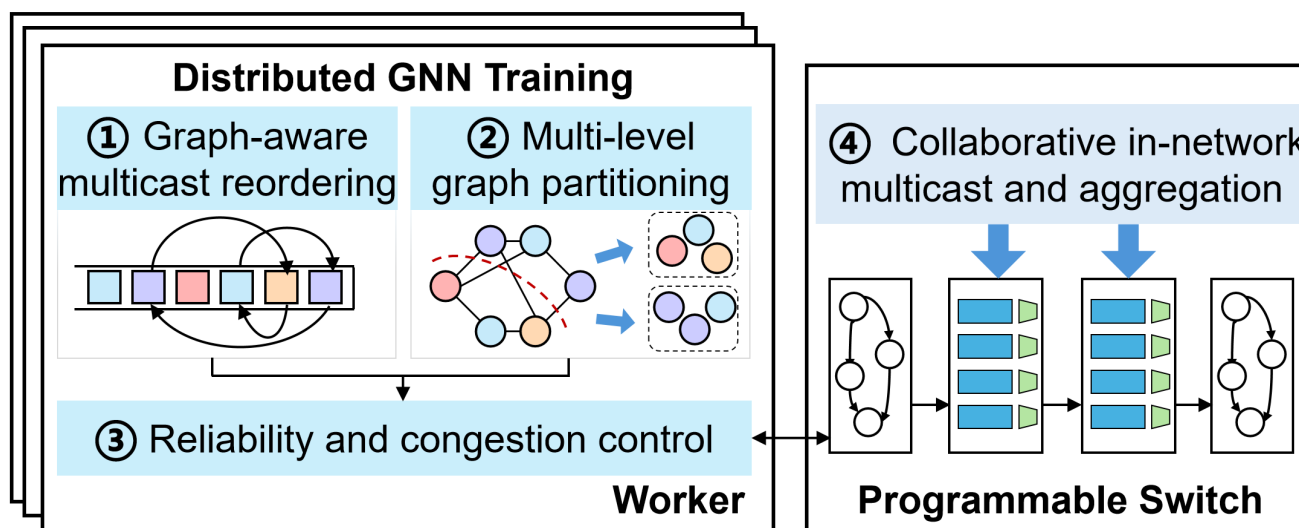
■ Deep Dive

- Graph-agnostic multicast order results in inefficient forwarding pipeline
- Huge and interdependent vertex features lead to aggregator overflow



SwitchGNN: Overview

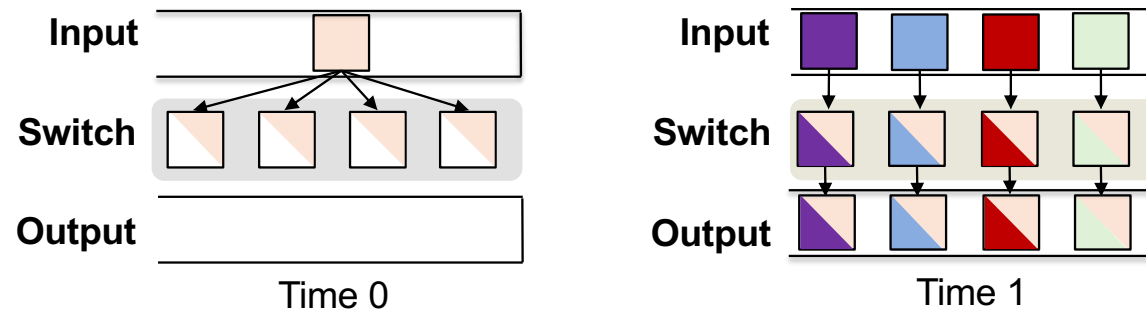
- Graph-aware multicast reordering optimizes the forwarding pipeline
- Multi-level graph partitioning mitigates aggregator overflow
- Reliability and congestion control avoids network congestion



SwitchGNN: Design Details

■ Graph-Aware Multicast Reordering

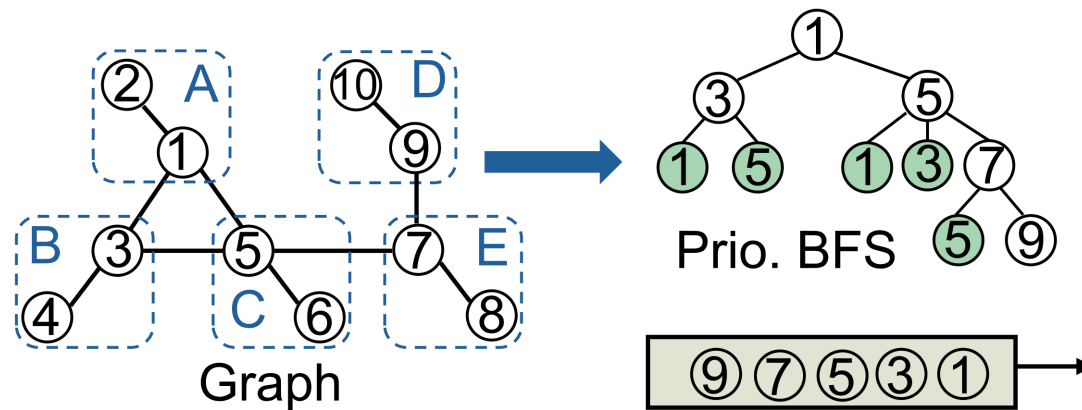
- Prioritizing transmission of high-outdegree vertex features improves forwarding efficiency



SwitchGNN: Design Details

■ Graph-Aware Multicast Reordering

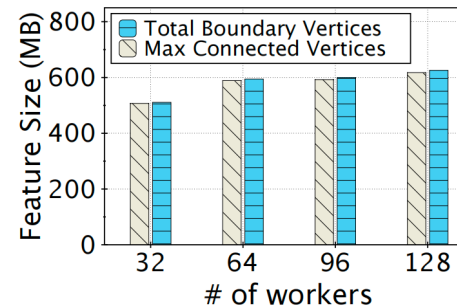
- Prioritizing transmission of high-outdegree vertex features improves forwarding efficiency
- Uploading features of neighboring vertices in close time enables faster aggregation at the switch



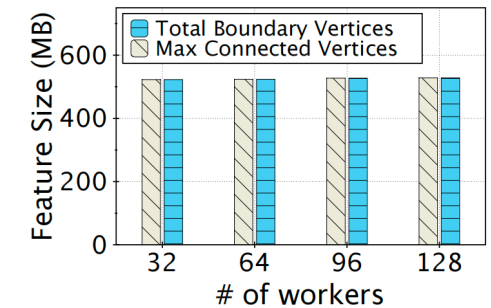
SwitchGNN: Design Details

■ Multi-level Graph Partitioning

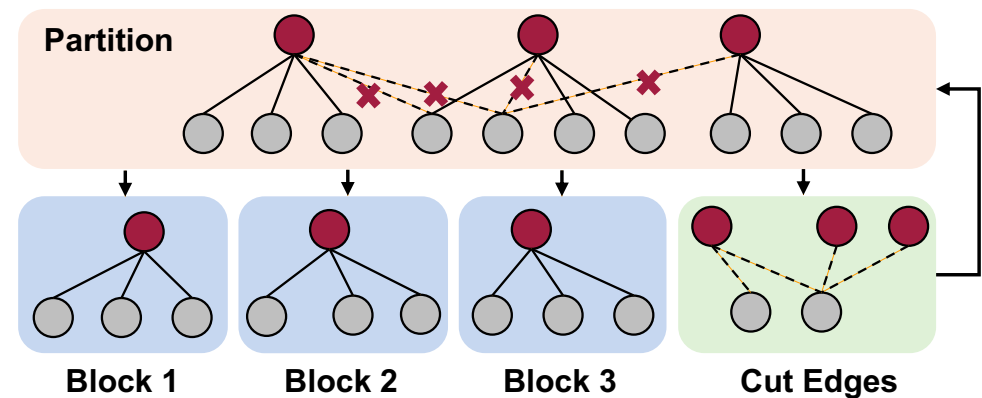
- Partitioning the graph (METIS) to fit switch memory and aggregating them in batches
- Subgraphs formed by cut edges should be still aggregated



(a) Ogbn-products



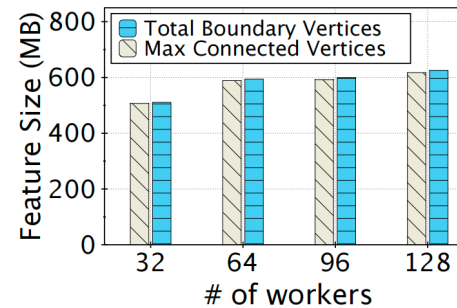
(b) Reddit



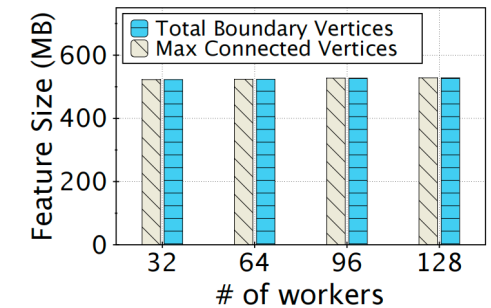
SwitchGNN: Design Details

■ Multi-level Graph Partitioning

- Partitioning the graph (METIS) to fit switch memory and aggregating them in batches
- Subgraphs formed by cut edges should be still aggregated



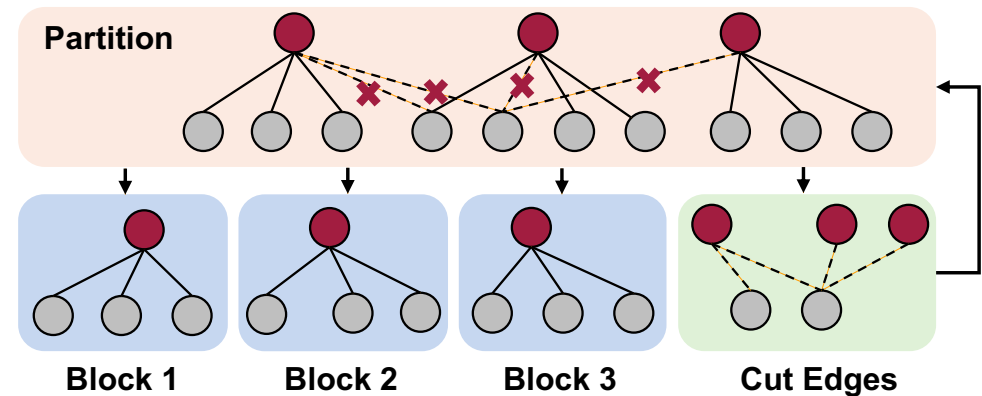
(a) Ogbn-products



(b) Reddit

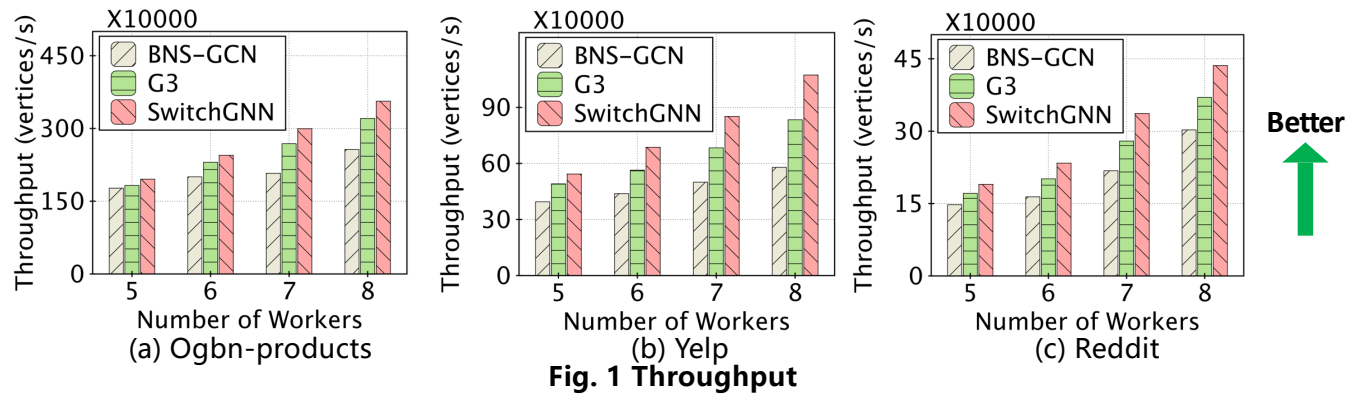
■ Reliability and Congestion Control

- Guaranteeing correct aggregation and avoid network congestion in all-to-all transmission

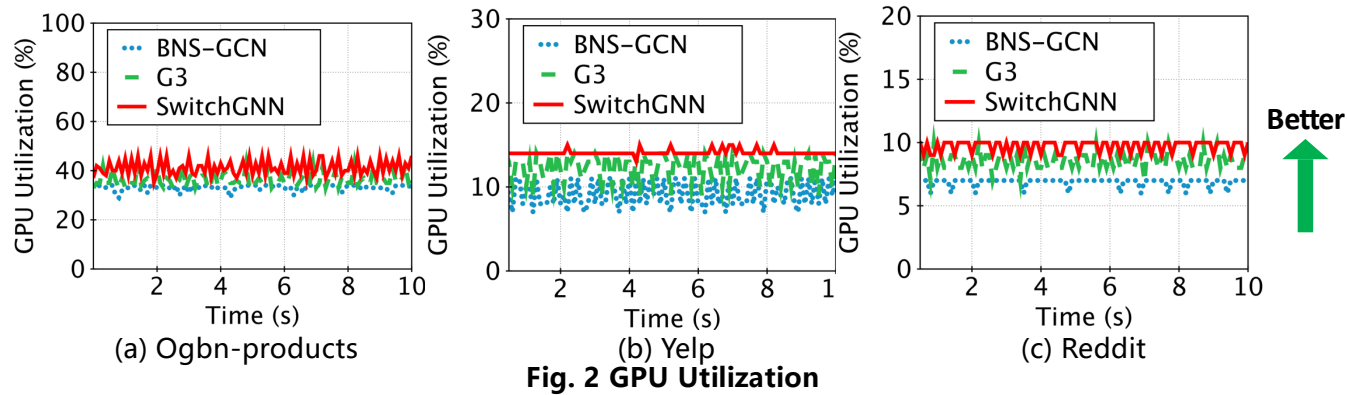


Evaluation

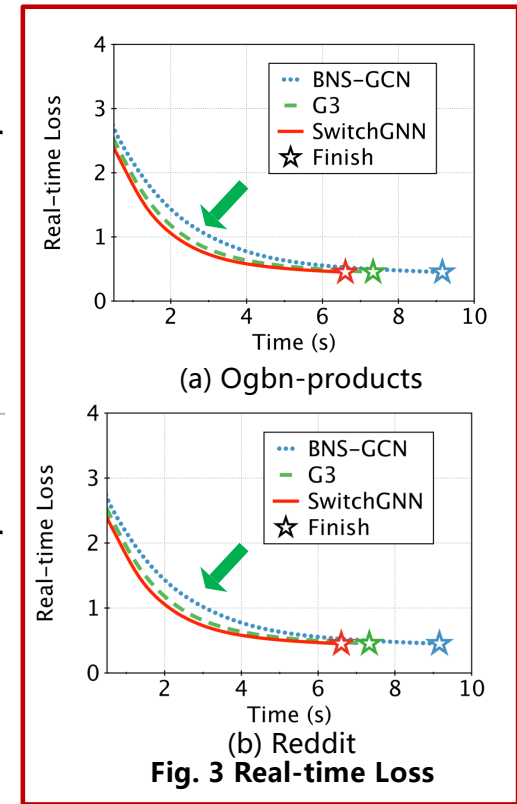
■ Testbed DPKD+P4; 8 servers; Star topology; GCN Model



Better
↑

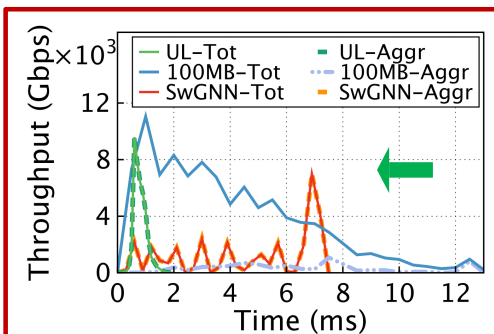


Better
↑

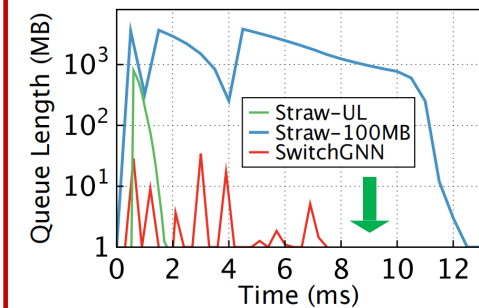


Evaluation

■ NS3 Simulation 128 workers; Leaf-spine topology; GCN model

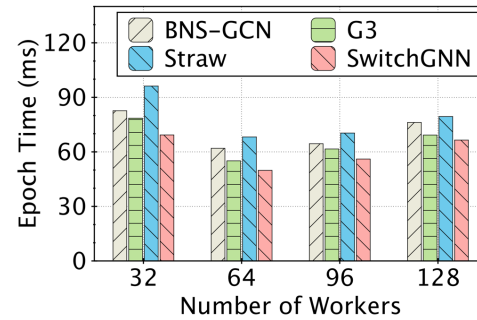


(a) Throughput

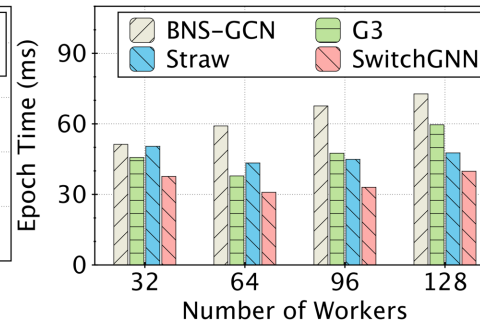


(b) Queue Length

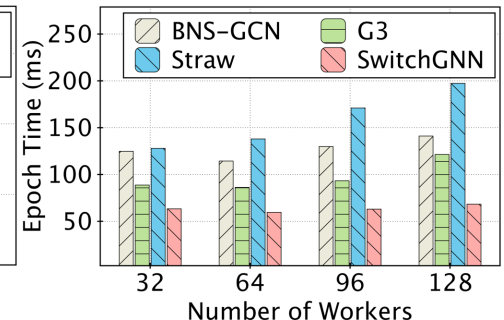
Fig. 4 Basic Performance



(a) Ognb-products

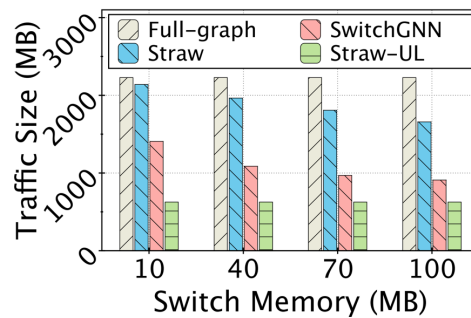


(b) Yelp

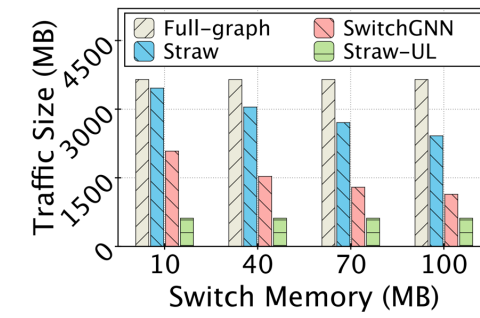


(c) Reddit

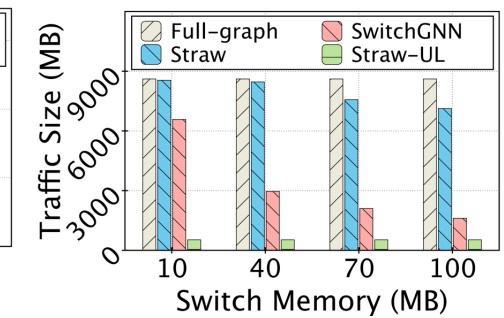
Fig. 5 Epoch Time



(a) Ognb-products



(b) Yelp

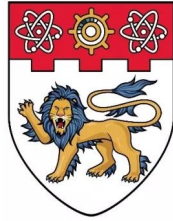


(c) Reddit

Fig. 6 Traffic Size

Summary

- SwitchGNN exploits in-network multicast and aggregation to accelerate GNN training in data centers
- Graph-aware multicast reordering to optimize the vertex transmission sequence for efficient in-network pipelining
- Multi-level graph partitioning to prevent aggregator overflow and minimize switch traffic
- **Key Takeaway: Harness ML workload characteristics to optimize both training efficiency and data center resource usage**



Thank You

chewei.tan@ntu.edu.sg

Accelerating distributed graph learning by using collaborative in-network multicast and aggregation

Zhaoyi Li^{1,3}, Jiawei Huang¹, Yijun Li¹, Jingling Liu¹, Junxue Zhang², Hui Li¹
Xiaojun Zhu¹, Shengwen Zhou¹, Jing Shao¹, Xiaojuan Lu¹, Qichen Su¹
Jianxin Wang¹, **Chee Wei Tan**³, Yong Cui⁴, Kai Chen²

¹Central South University, ²Hong Kong University of Science and Technology

³Nanyang Technological University, ⁴Tsinghua University