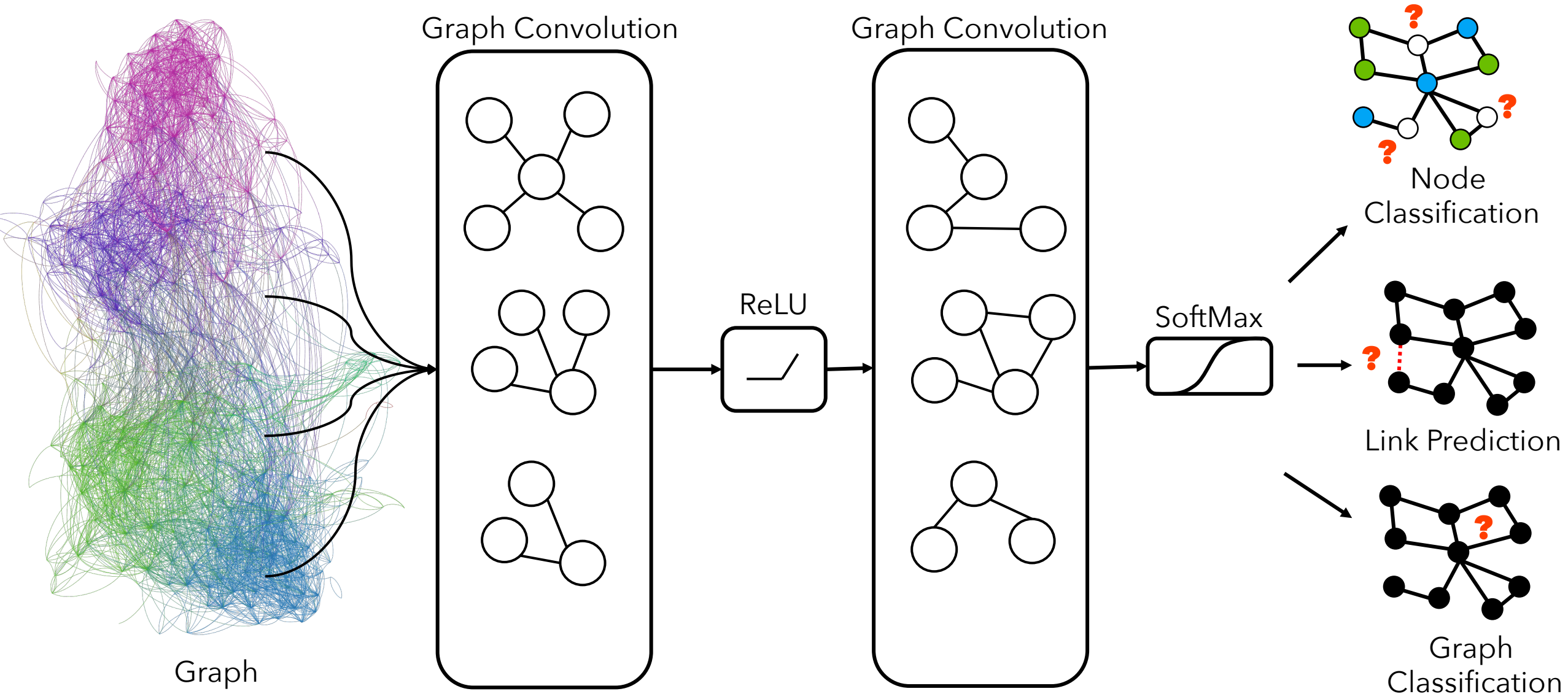


P³: Distributed Deep Graph Learning at Scale

Swapnil Gandhi, Anand Iyer
Microsoft Research

OSDI 2021

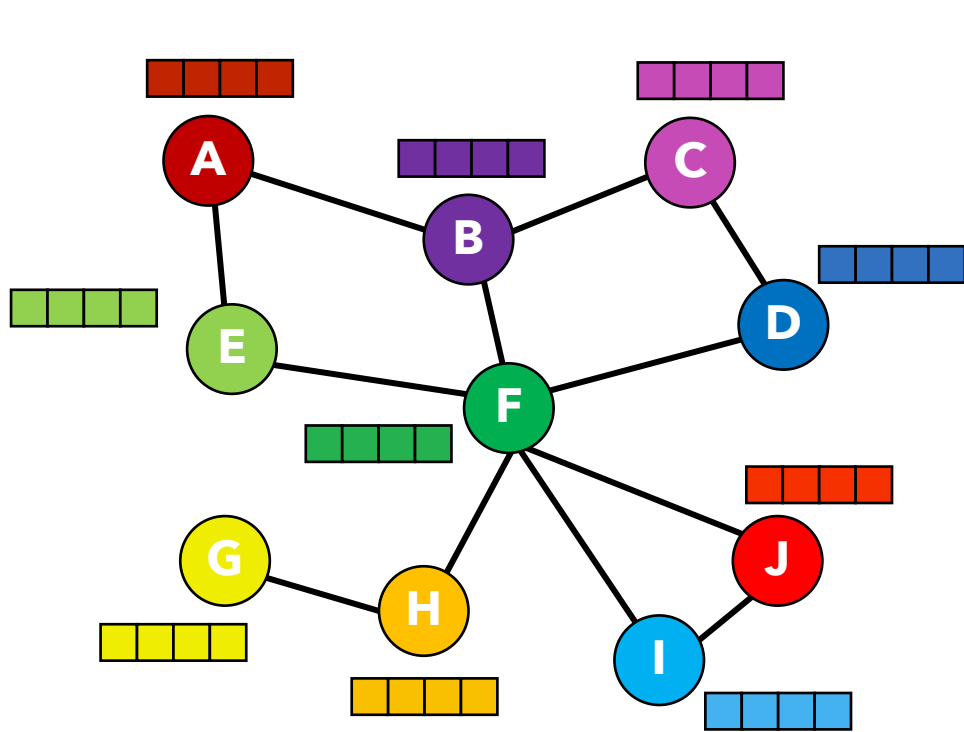
Graph Neural Networks



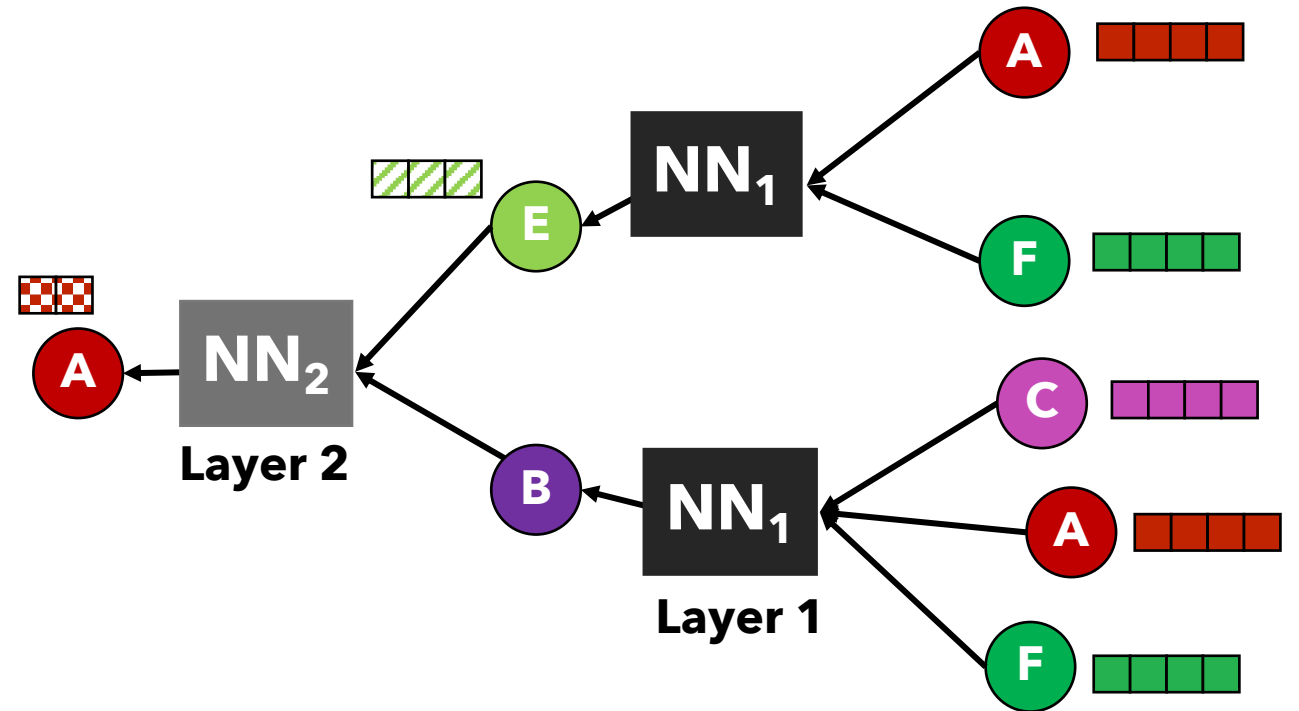
Graph Neural Networks

Graph Structure: What (to propagate?)

Neural Network: How (information is transformed)



Input Graph



2-hop computation graph of A

GNN Training

Large Graphs

Millions of nodes,
billions of edges

**Hundreds or
thousands** of features

New Models

Several proposals: GCN,
GAT, GIN, ...

More **sophisticated**,
complex architectures

Significant interest in **distributed** GNN training

**Distributed Graph
Neural Networks**

=

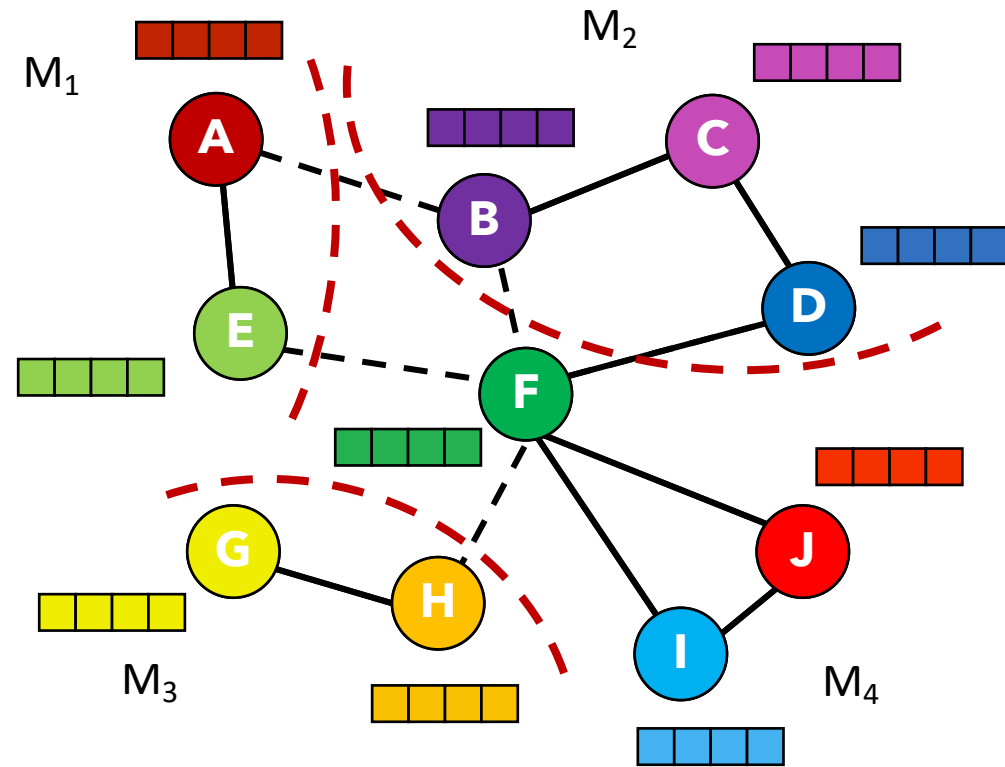
**Distributed Graph Processing
Techniques**

+

**Distributed Neural Network
Techniques**

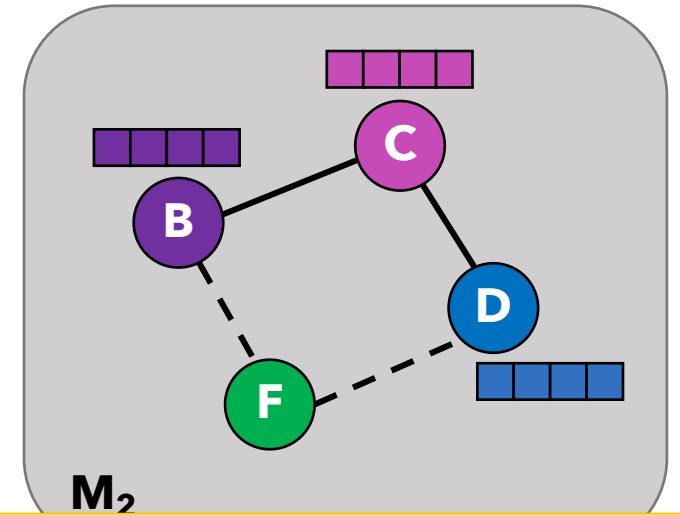
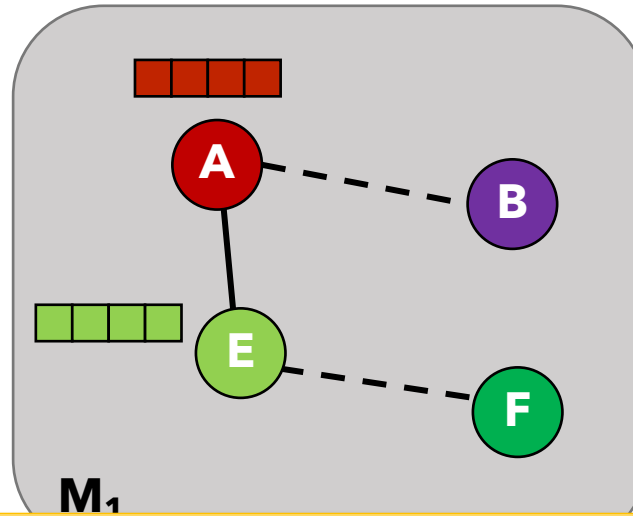
Distributed Graph Processing Techniques

Graph Partitioning

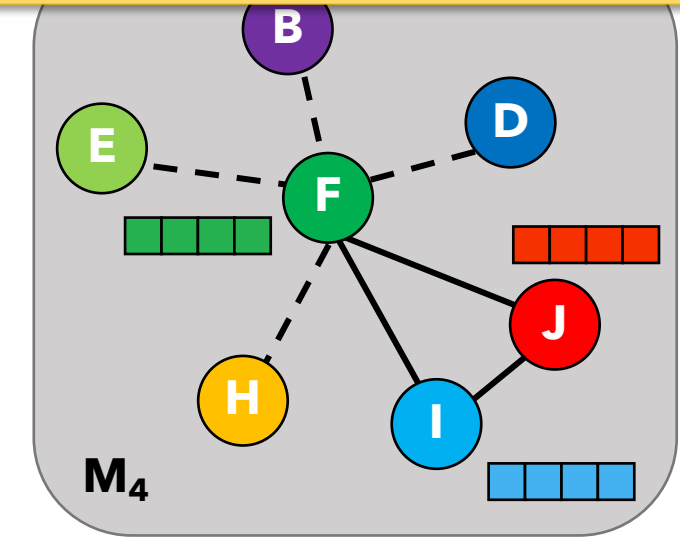
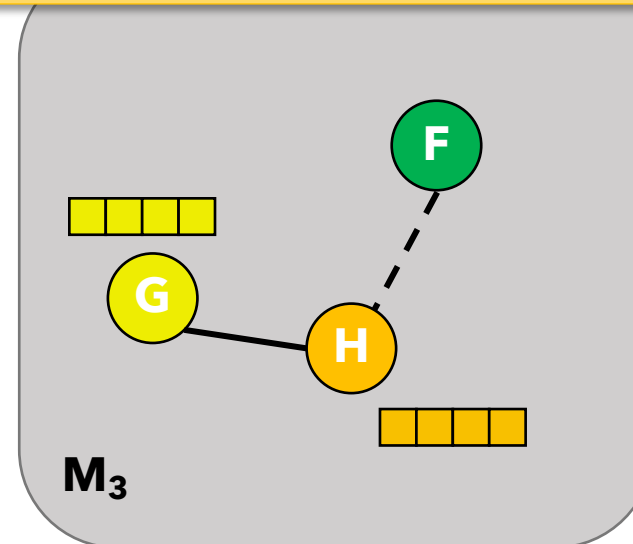


Distributed Graph Processing Techniques

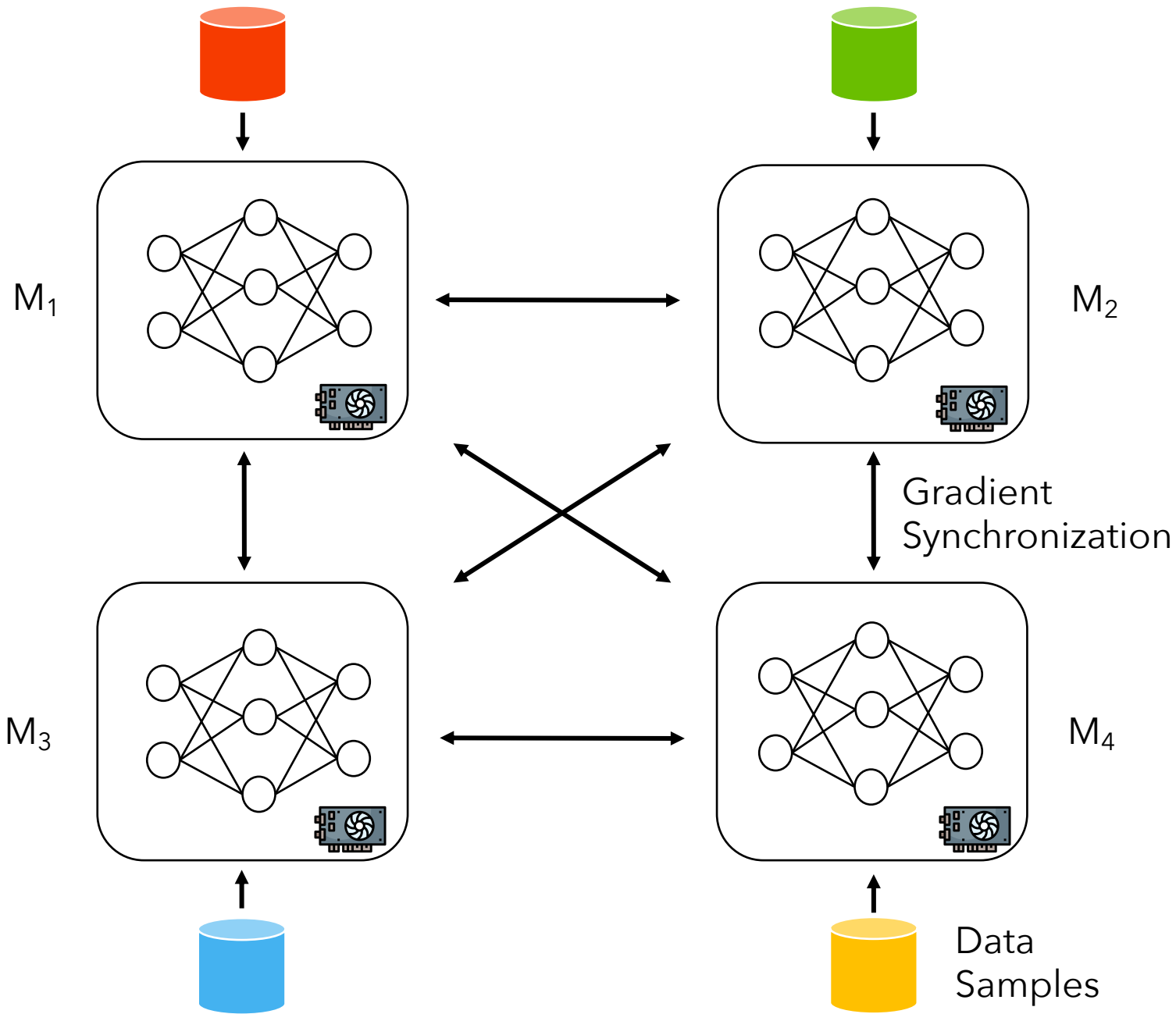
Graph Partitioning



Edge-cut, vertex-cut, hybrid, ...



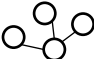


— Local Edge - - - Cut Edge 6

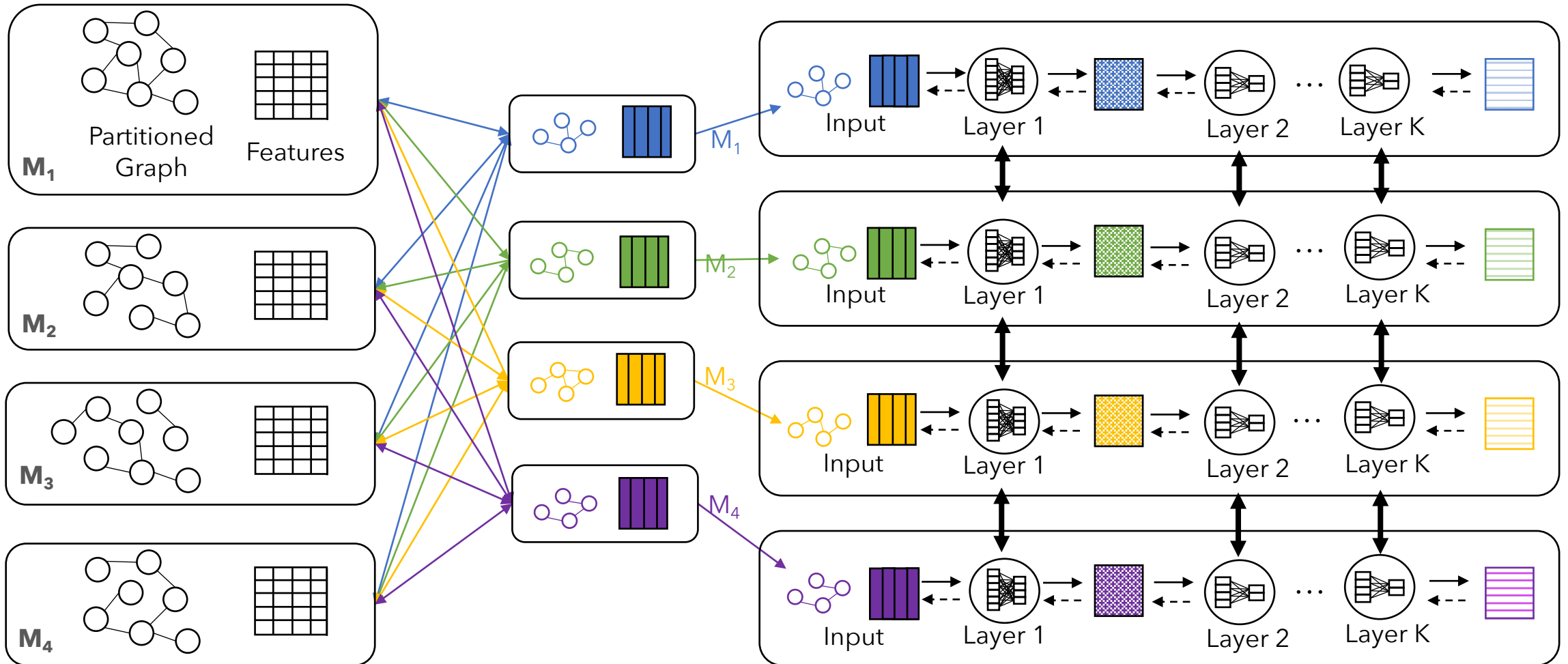


Data-parallel Training

Distributed Neural Network Processing Techniques

Distributed Graph Neural Networks

 GNN Computational Graph  Input Features \rightarrow Forward Pass \leftarrow Backward Pass \leftrightarrow Gradient Sync  Embedding



Network overhead **dominates** epoch time,
rendering **GPUs underutilized ~80%** of the time

Partitioning is **ineffective**, and in many cases
counterproductive

P³ proposes **push-pull parallelism**, a new
technique for distributed GNN training that
effectively **eliminates** these overheads

P³: Pipelined Push Pull

Feature movements cause dominant network traffic

Graph structure can be **compactly** represented

Reduce data communicated by **avoiding** feature movement

Existing systems consider graph & features **indivisible**

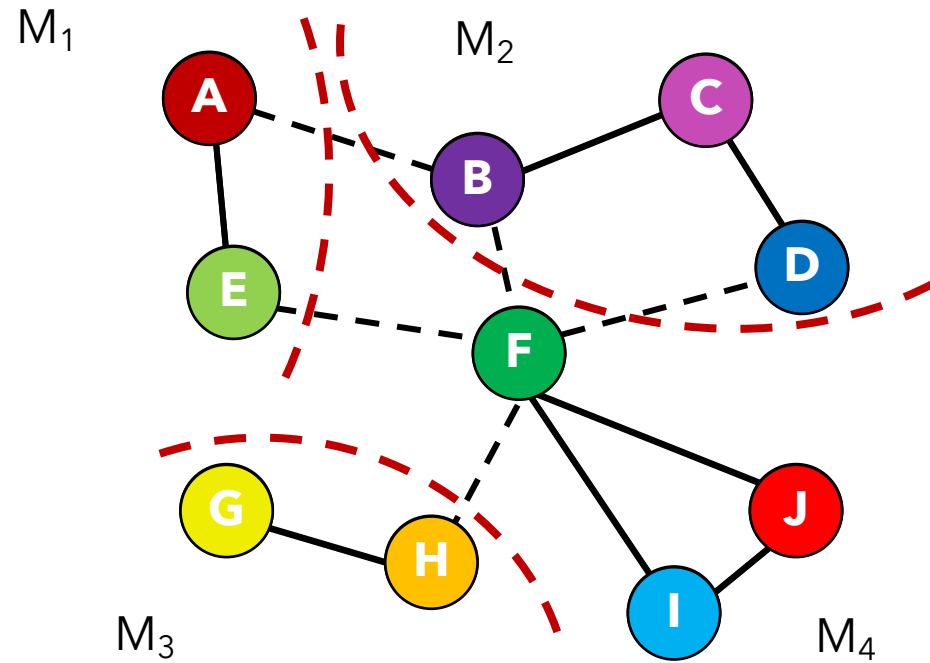
Partition them as **single-entity**

Independent Partitioning of Graph & Features

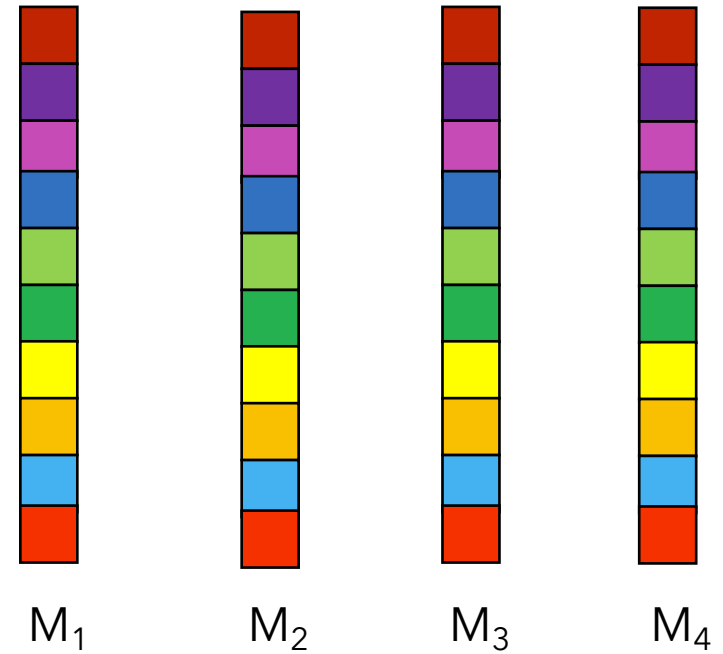
Independent Partitioning

Graph Structure: Partitioned using random hash func

Features: Partitioned along feature dimension



Graph Structure

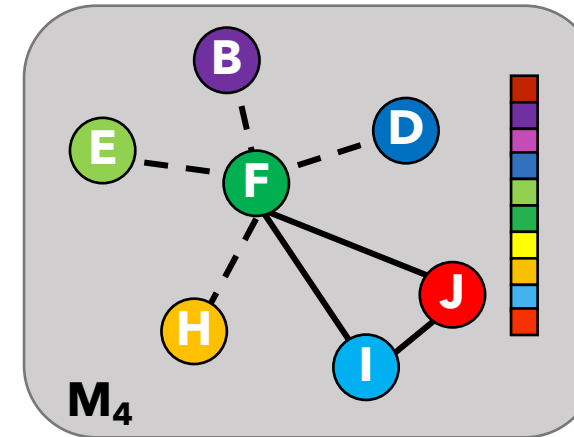
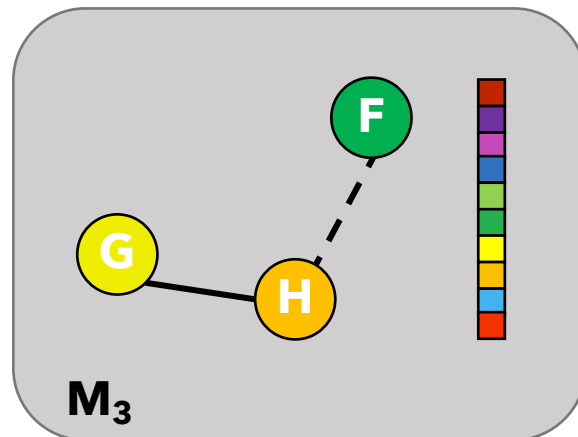
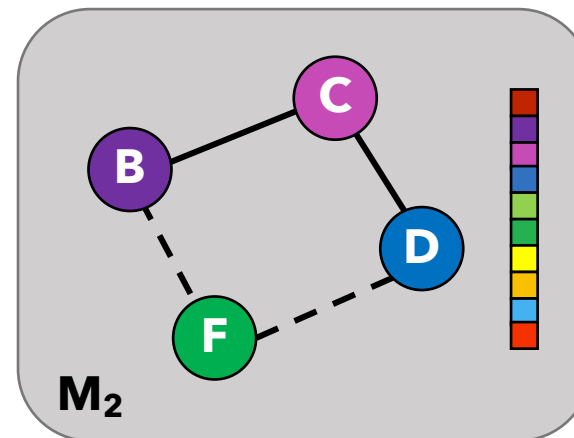
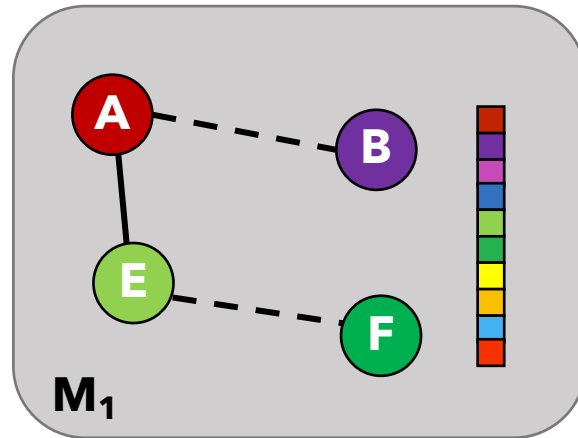


Features

Independent Partitioning

Graph Structure: Partitioned using random hash func

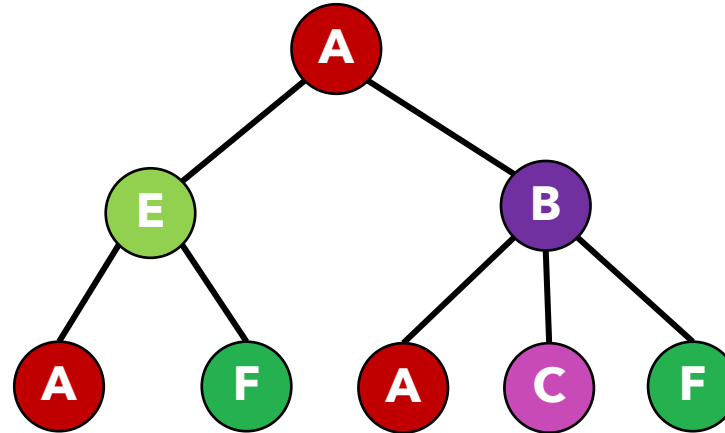
Features: Partitioned along feature dimension



Computation Graph Creation

→ Local → Remote

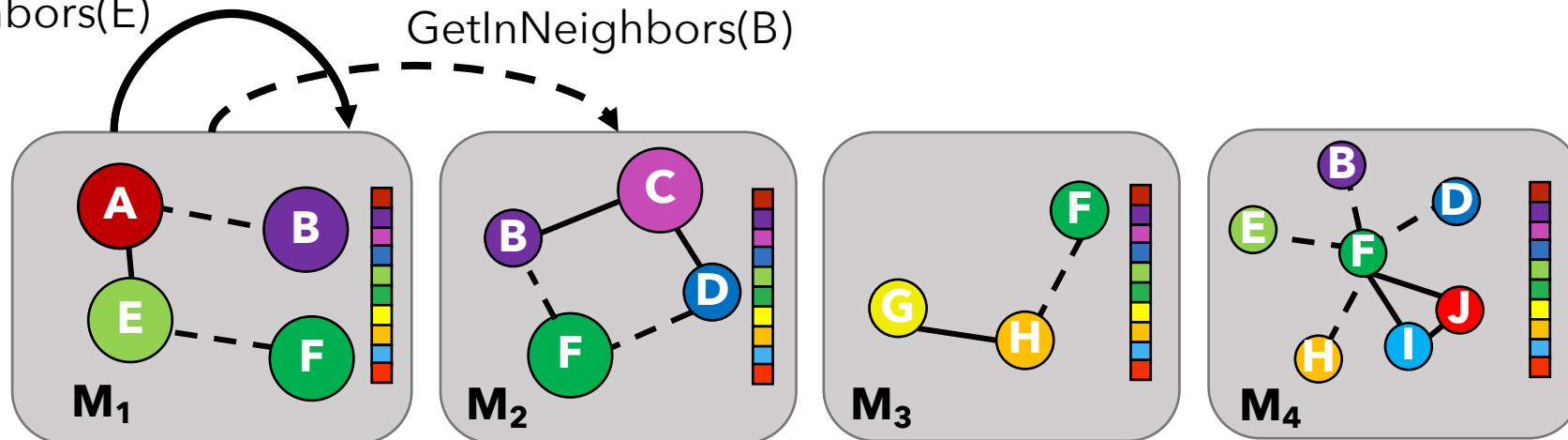
Structure moved over network



Features **not** moved over network

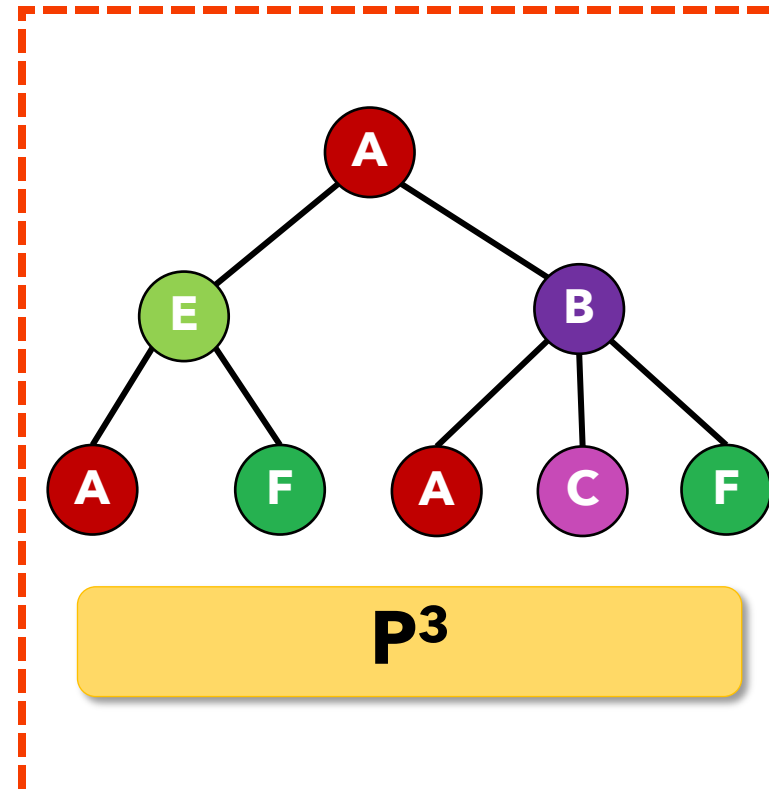
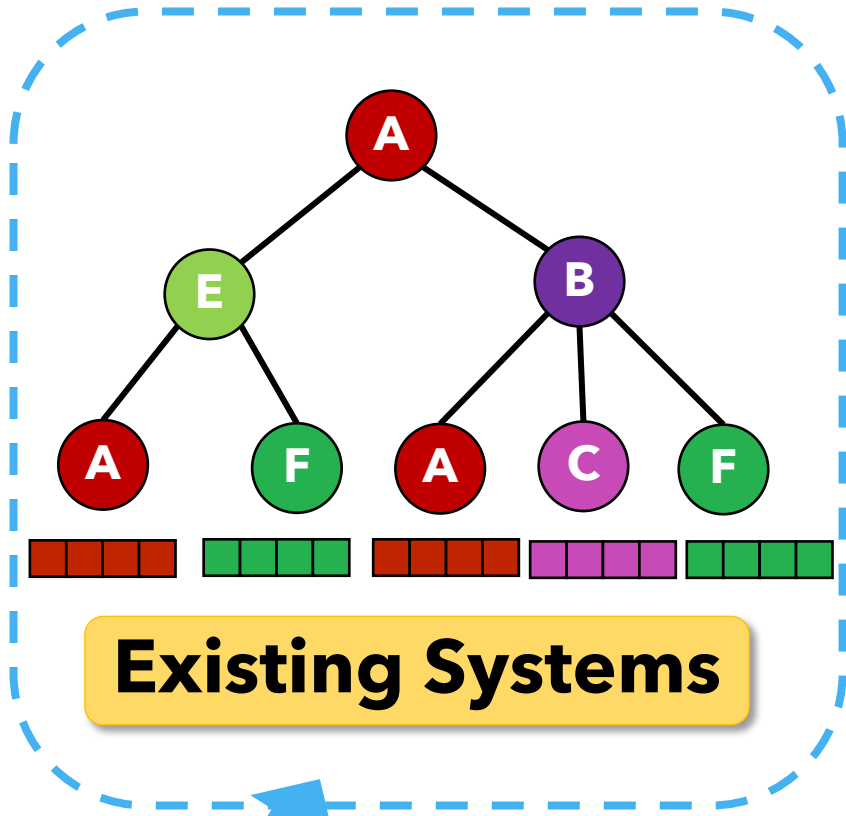
GetInNeighbors(E)

GetInNeighbors(B)



Results in **significant reduction** in data communication

Hybrid Parallelism

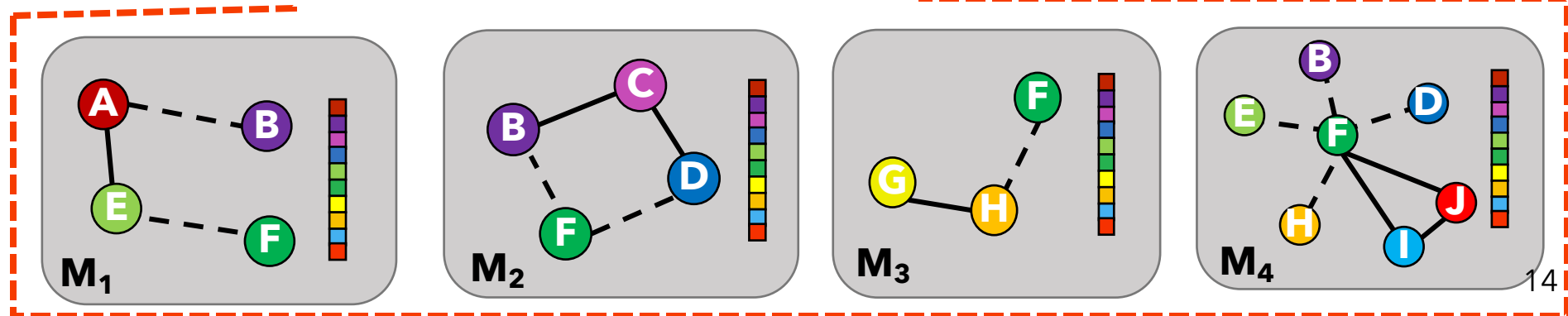


Cannot use data parallelism

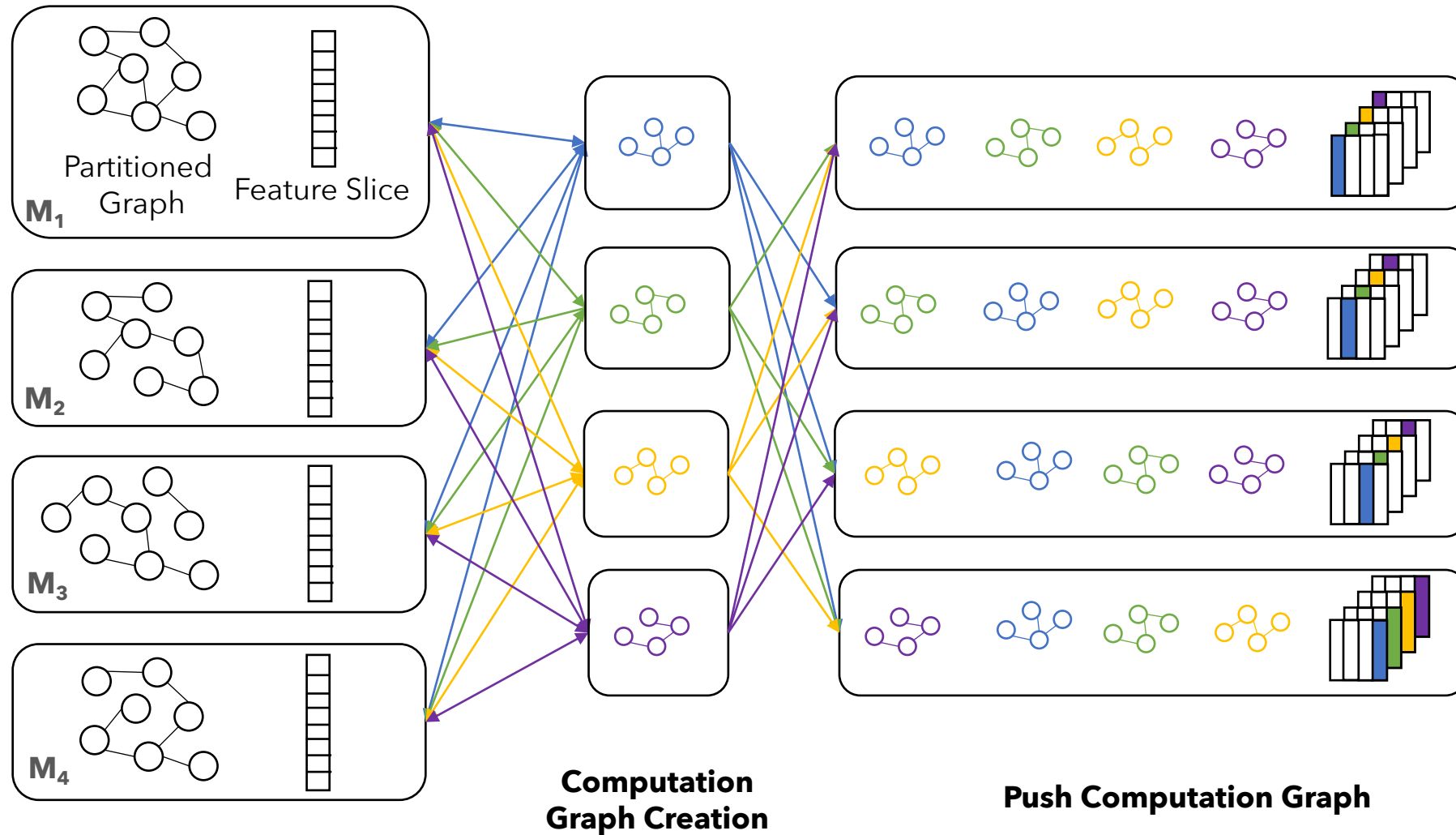
Model parallelism incurs overheads

P³ combines data and model parallelism

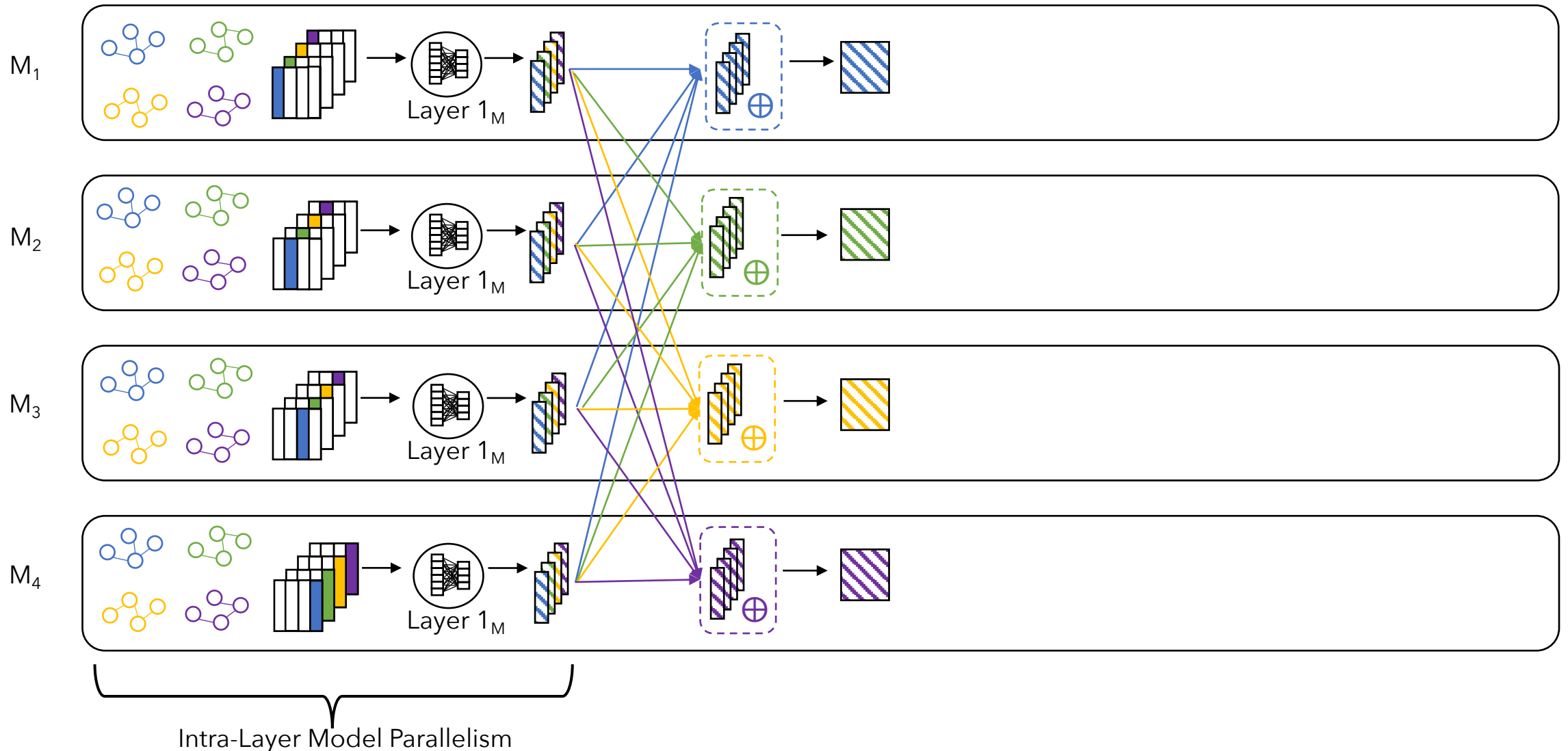
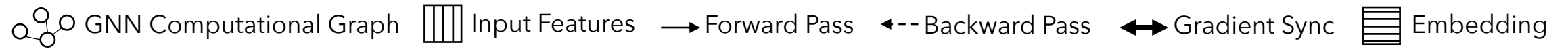
Use Data Parallelism



P³: Distributed Graph Neural Networks



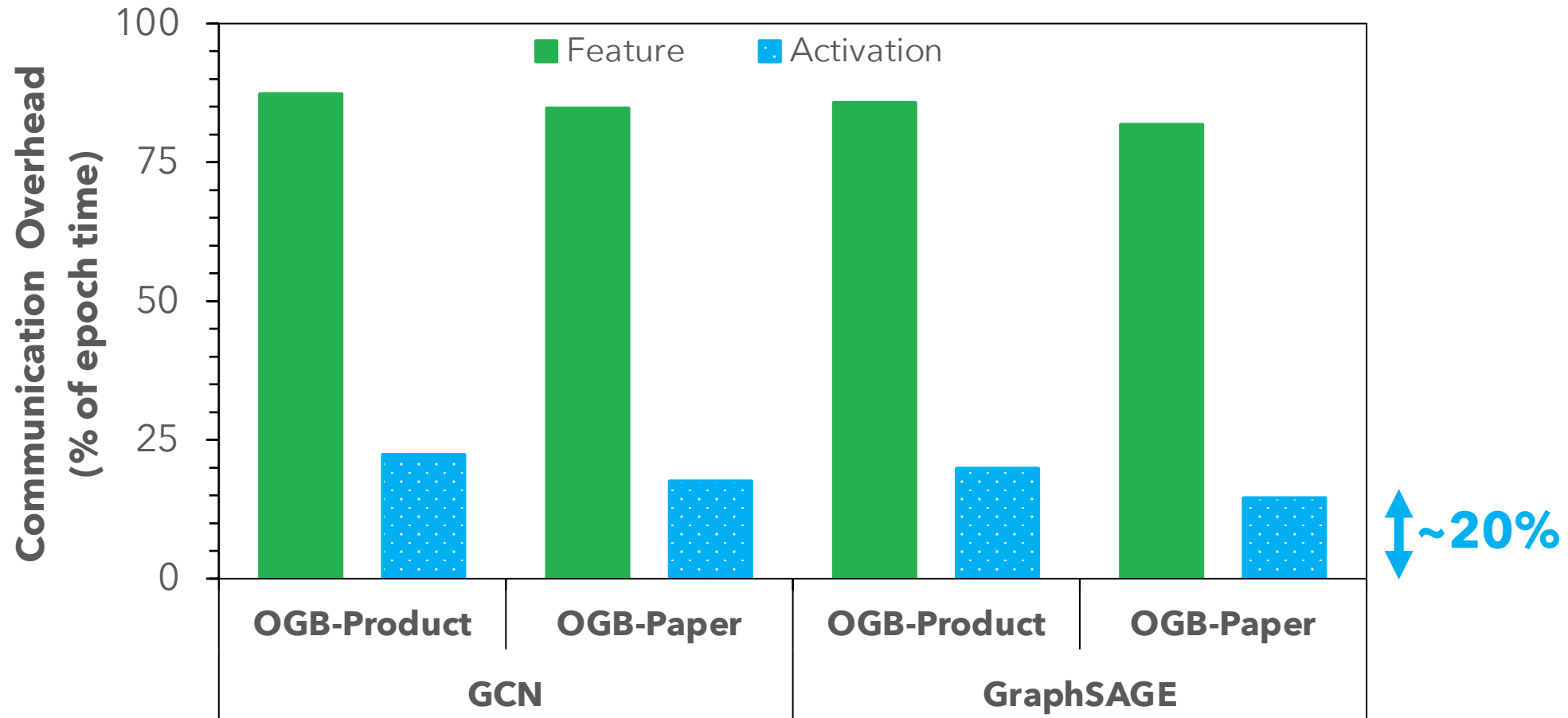
P³: Distributed Graph Neural Networks



GNNs typically use **small** hidden dimensions

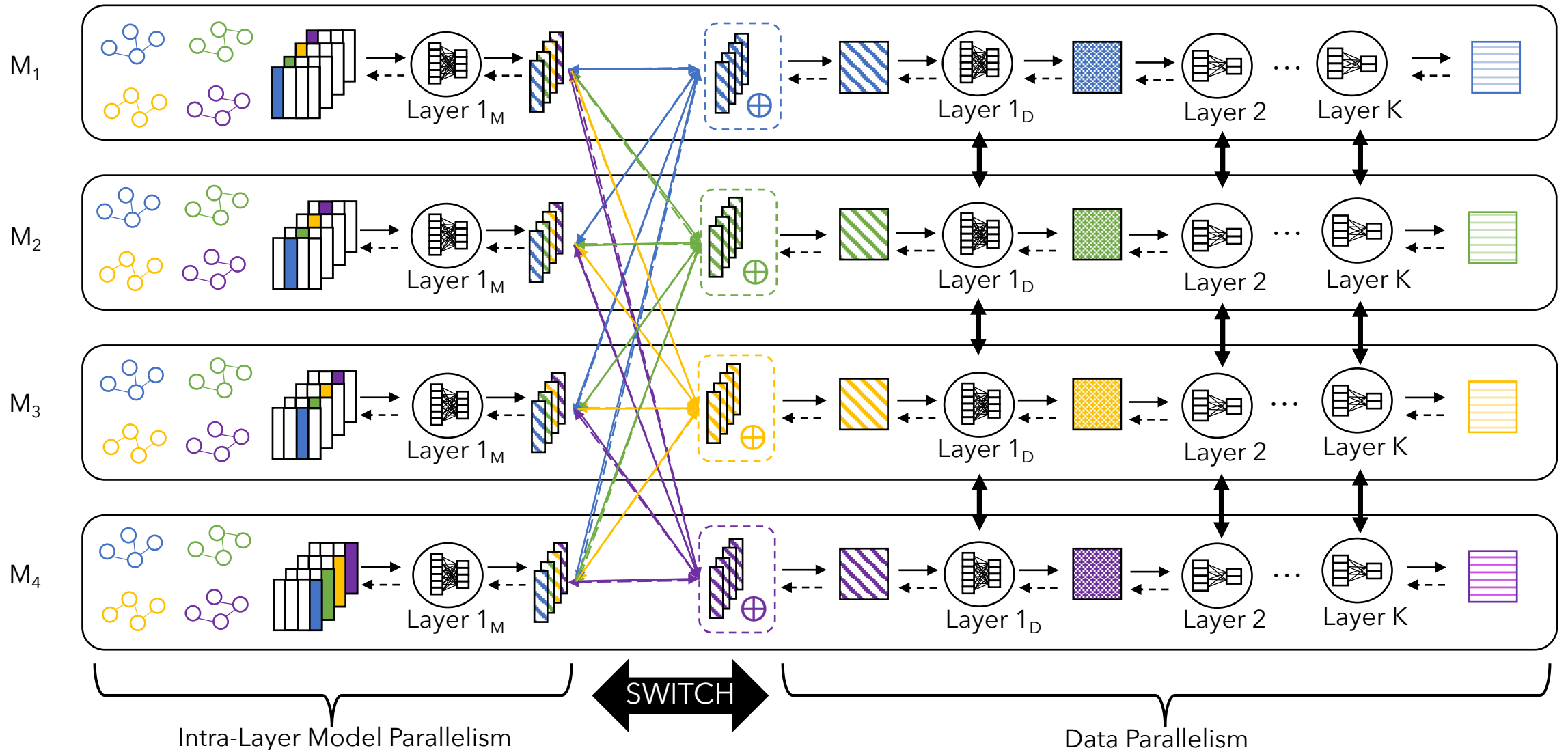
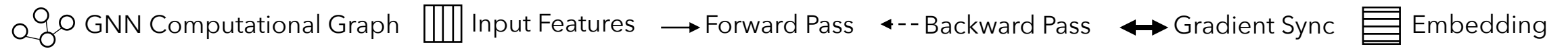


small intermediate activations



Enables **faster** training across range of GNN models and datasets

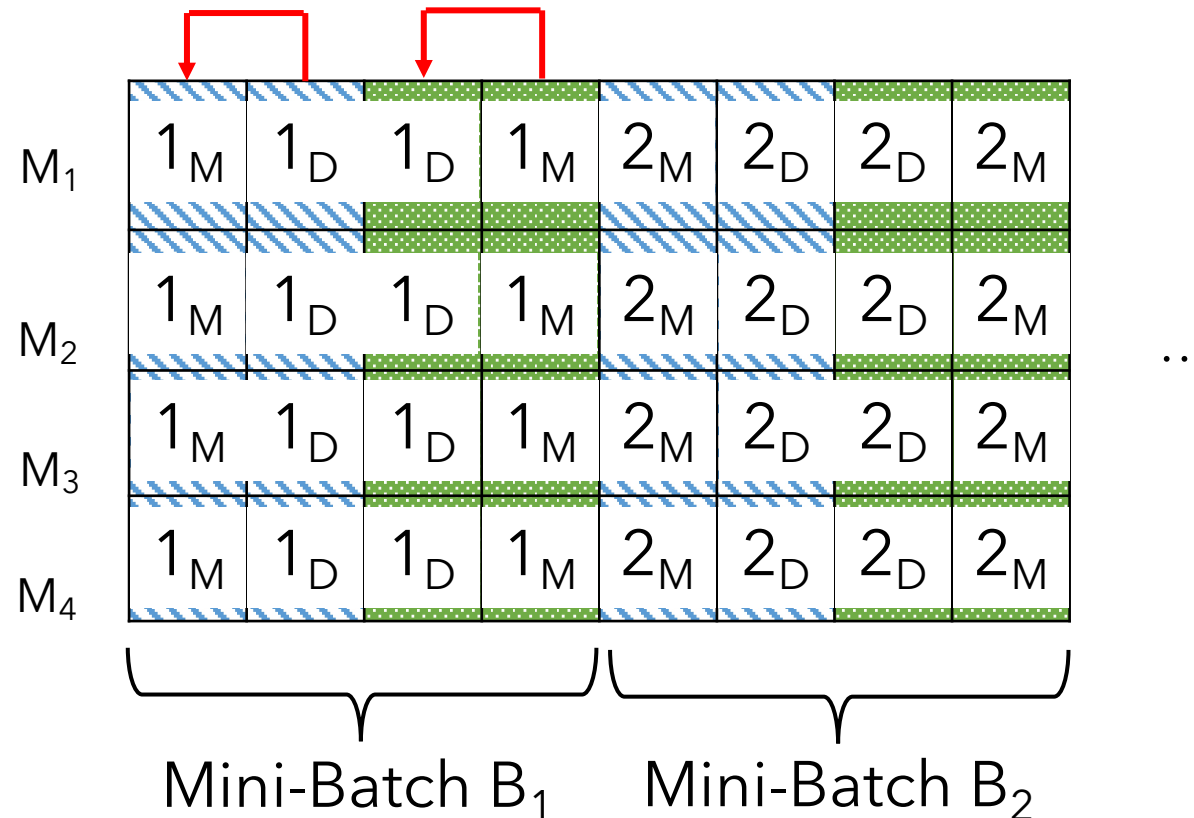
P³: Distributed Graph Neural Networks



P³: GNN Training

Hybrid Parallelism requires **communication** in **forward** and **backward** pass

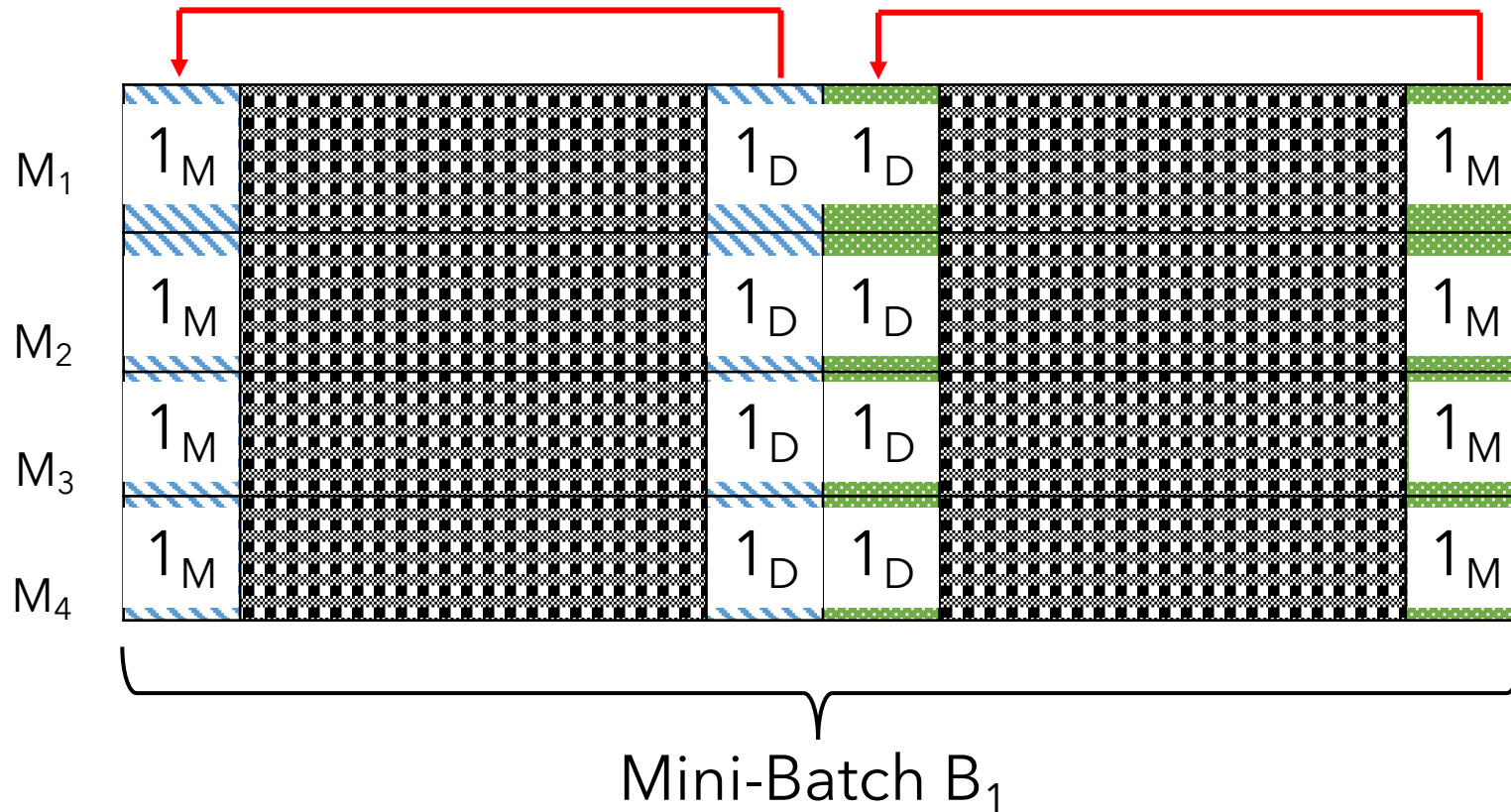
 Forward Pass  Backward Pass  Data Dependency



P³: GNN Training

Communication results in **computation stall**

 Forward Pass  Backward Pass  Stall  Data Dependency

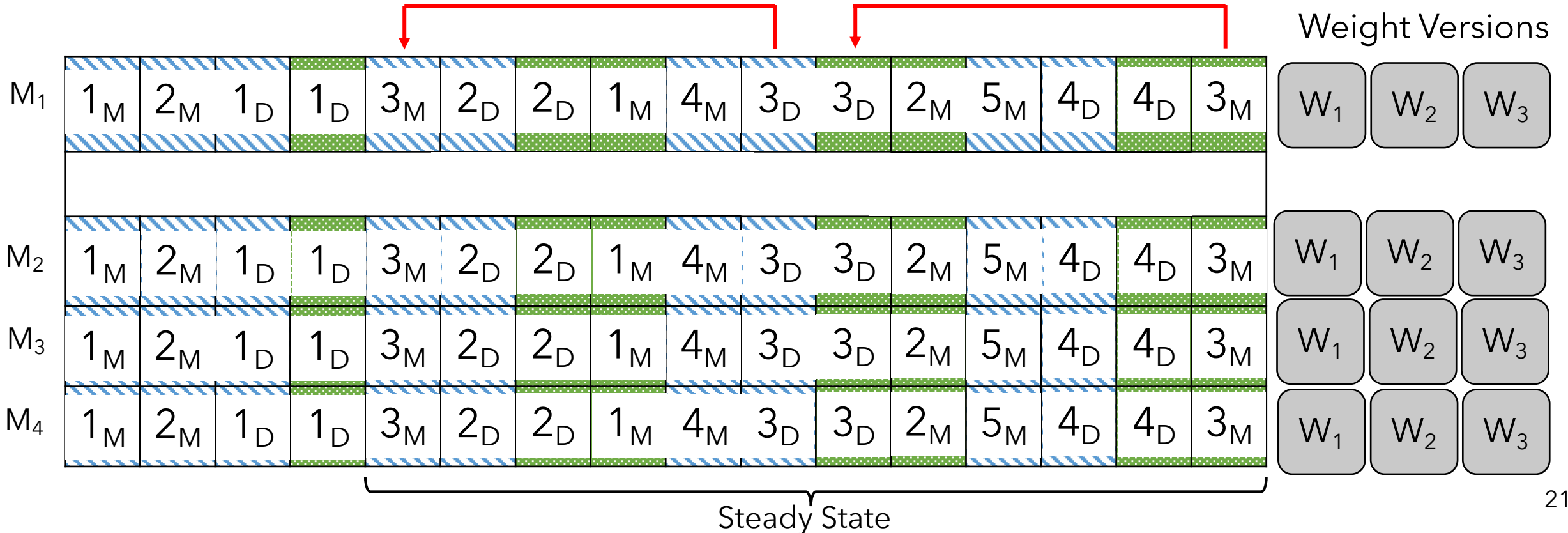


P3: Pipelining

Overlaps **computation** with **communication**

Improves performance by up to **50%**

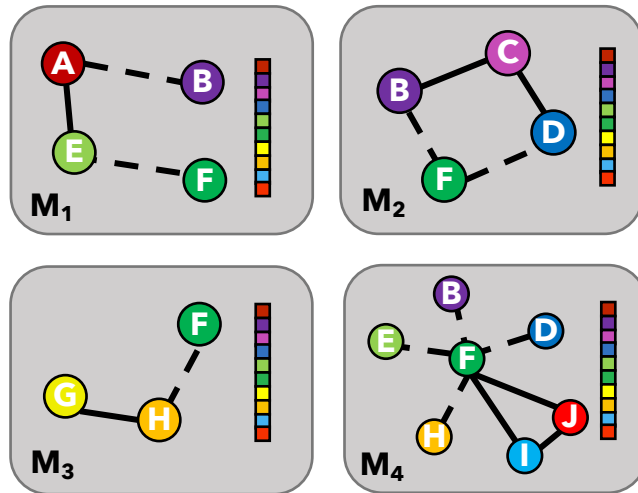
 Forward Pass
  Backward Pass
  Stall
  Data Dependency



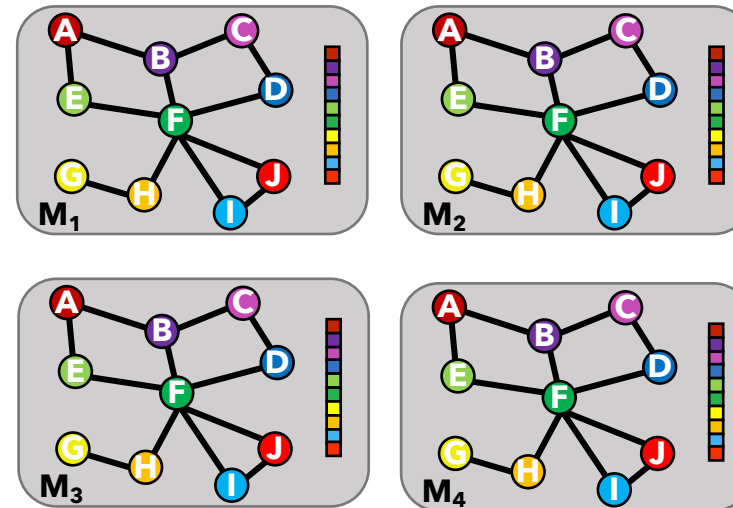
P³: Caching

Cache **graph structure** and/or **features**

Improves performance by up to **1.7x**



Structure : Partitioned
Feature : Partitioned



Structure : Cached
Feature : Partitioned

P³: API

P³ exposes a **simple API** for developers

partition()

Independently partition graph and features, and cache if possible

scatter()

Generate message vector

gather()

Aggregate message vector

transform()

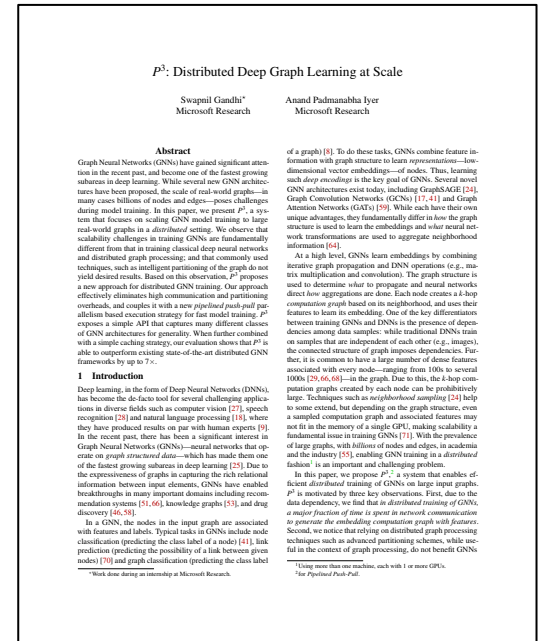
Compute partial activation

sync()

Accumulate partial activation

apply()

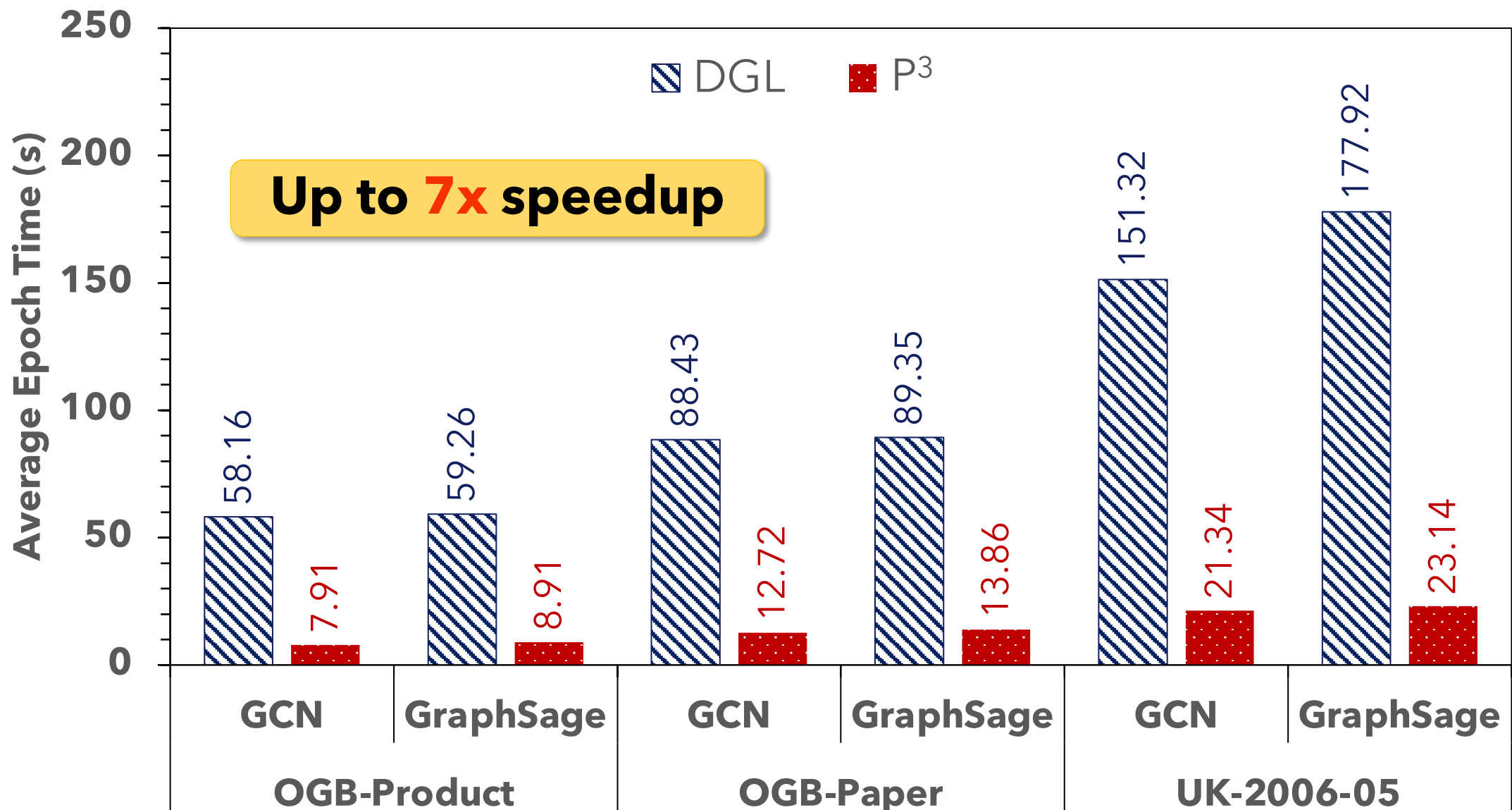
Compute output activation



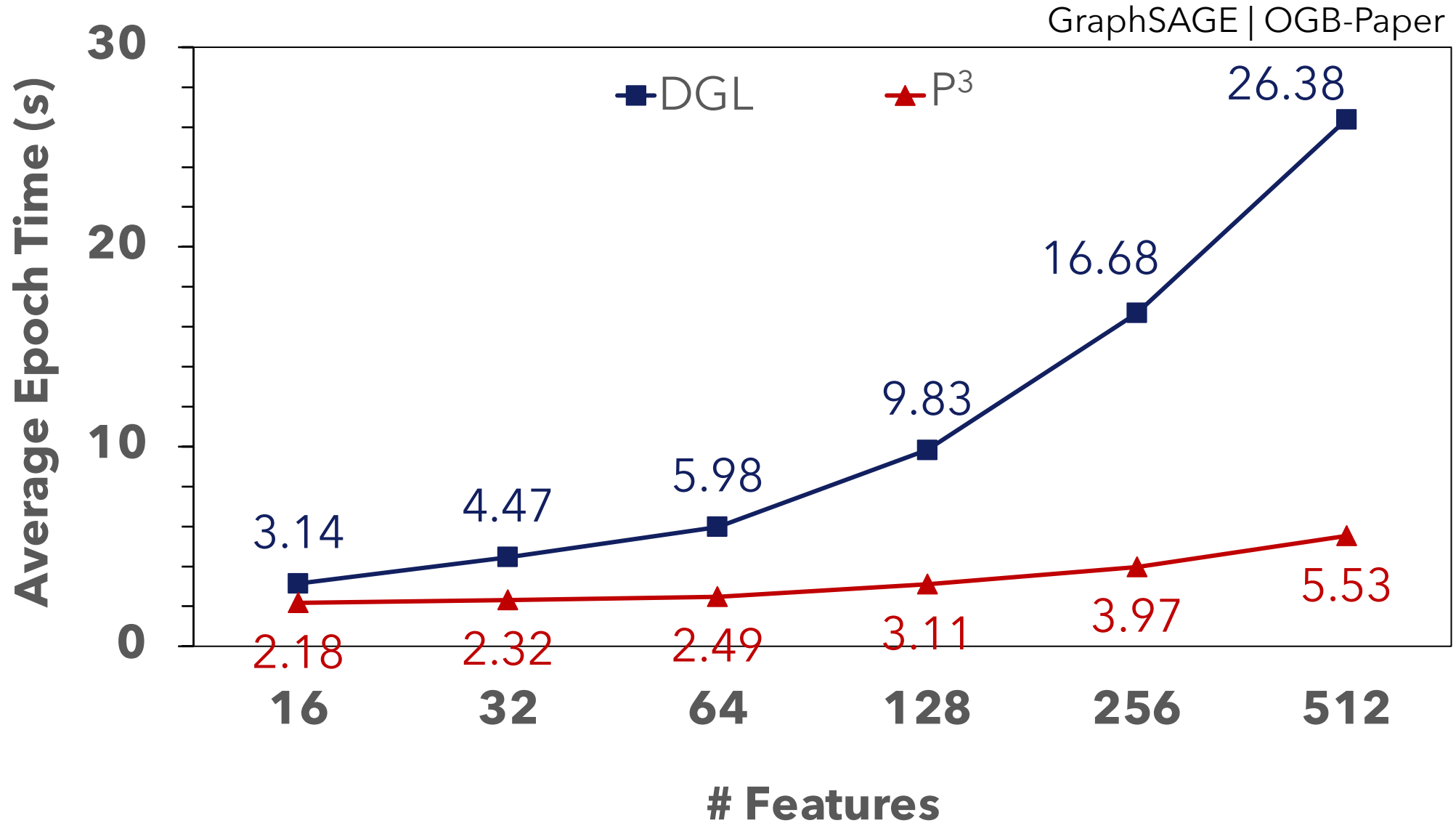
Implementation & **Evaluation**

- Implemented on Deep Graph Library (DGL) v0.5
 - Uses PyTorch v1.6
- Evaluated using 16 NVIDIA Tesla P100 GPUs
 - **OGN-Product:** 123M edges, $|F|=100$
 - **OGN-Paper:** 1.6B edges, $|F|=128$
 - **UK-2006-05:** 2.9B edges, $|F|=256$
- Several GNN architectures
 - GCN [NeurIPS '16]
 - GraphSAGE [NeurIPS '17]

P³ Performance

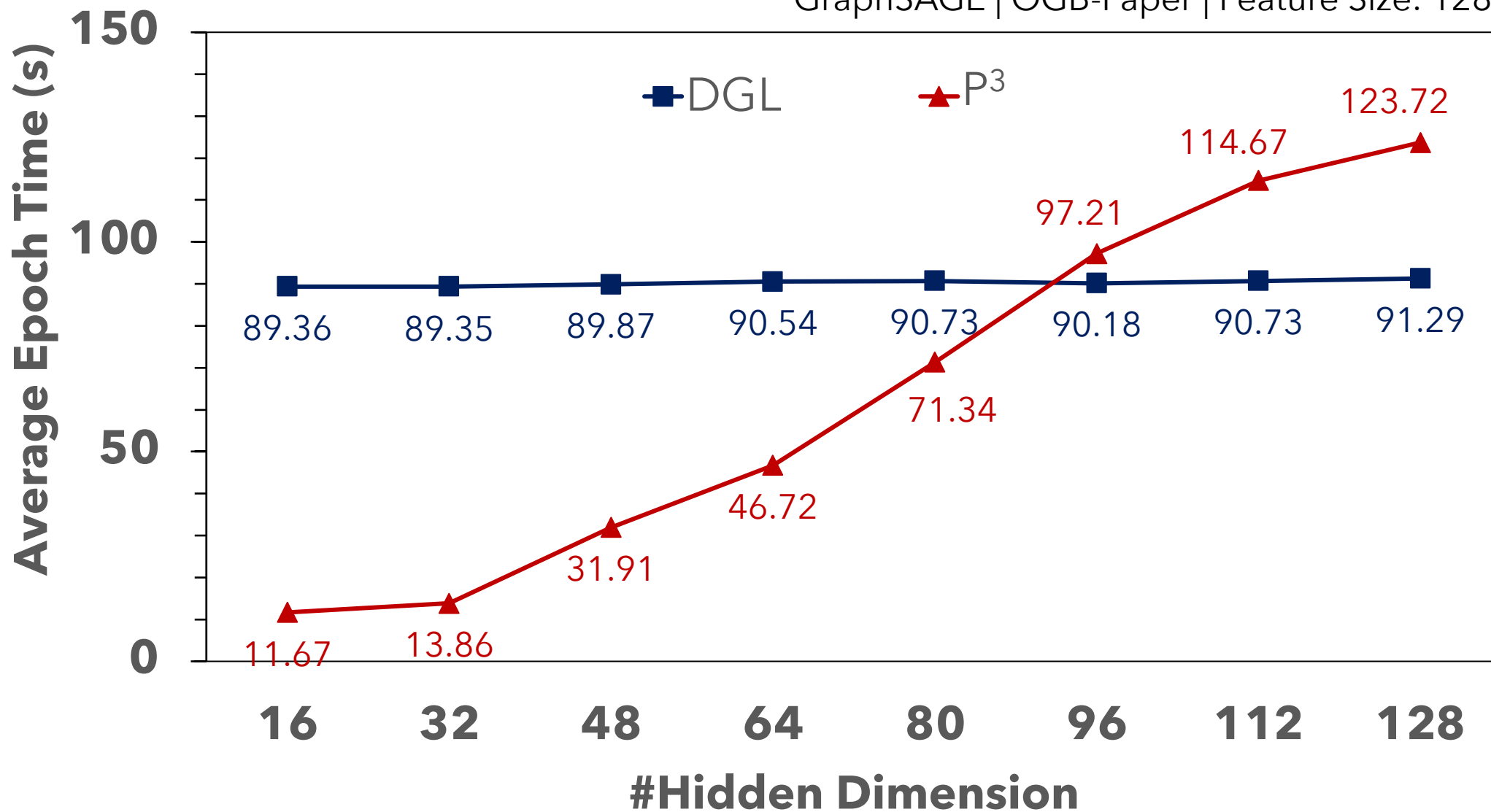


P³ Scaling



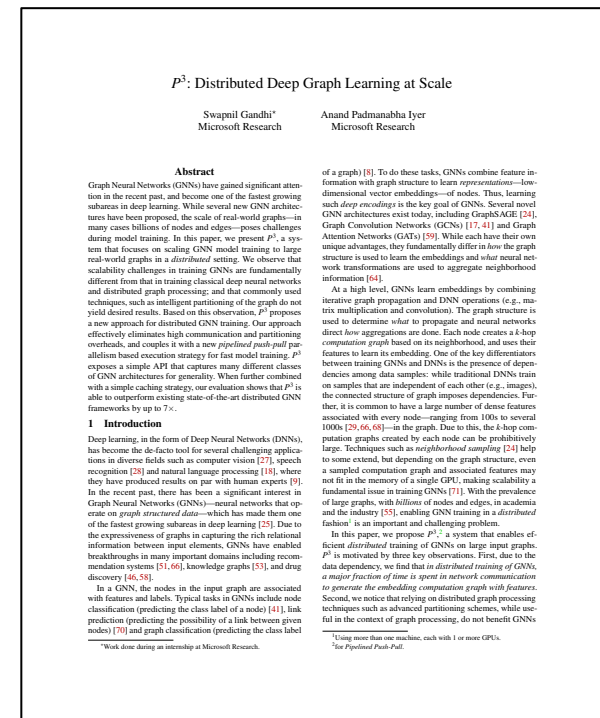
P³ Shortcomings

GraphSAGE | OGB-Paper | Feature Size: 128



More evaluation in the paper

- GNN Models: GAT [ICLR'18], SGCN [ICML'19]
- Larger Datasets: UK-Union ($|E|=5.5B$), Facebook ($|E|=10B$)
- Study impact of
 - Sampling
 - Partitioning Strategies
 - Number of Layers
 - Pipelining
 - Caching
- Scaling Characteristics
- Comparison with ROC [MLSys'20]



Takeaway

- Distributed training of graph neural networks increasingly important
 - Frameworks = Graph Processing + DNN Training
 - Incur high network communication and partitioning overhead
- P³ eliminates the overheads with distributed GNN training
 - **Independent partitioning** of graph structure and features
 - **Hybrid parallelism** combined with **pipelining** and **caching**
 - Simple **API** for users

Thank you!
<https://swapnilgandhi.com>