

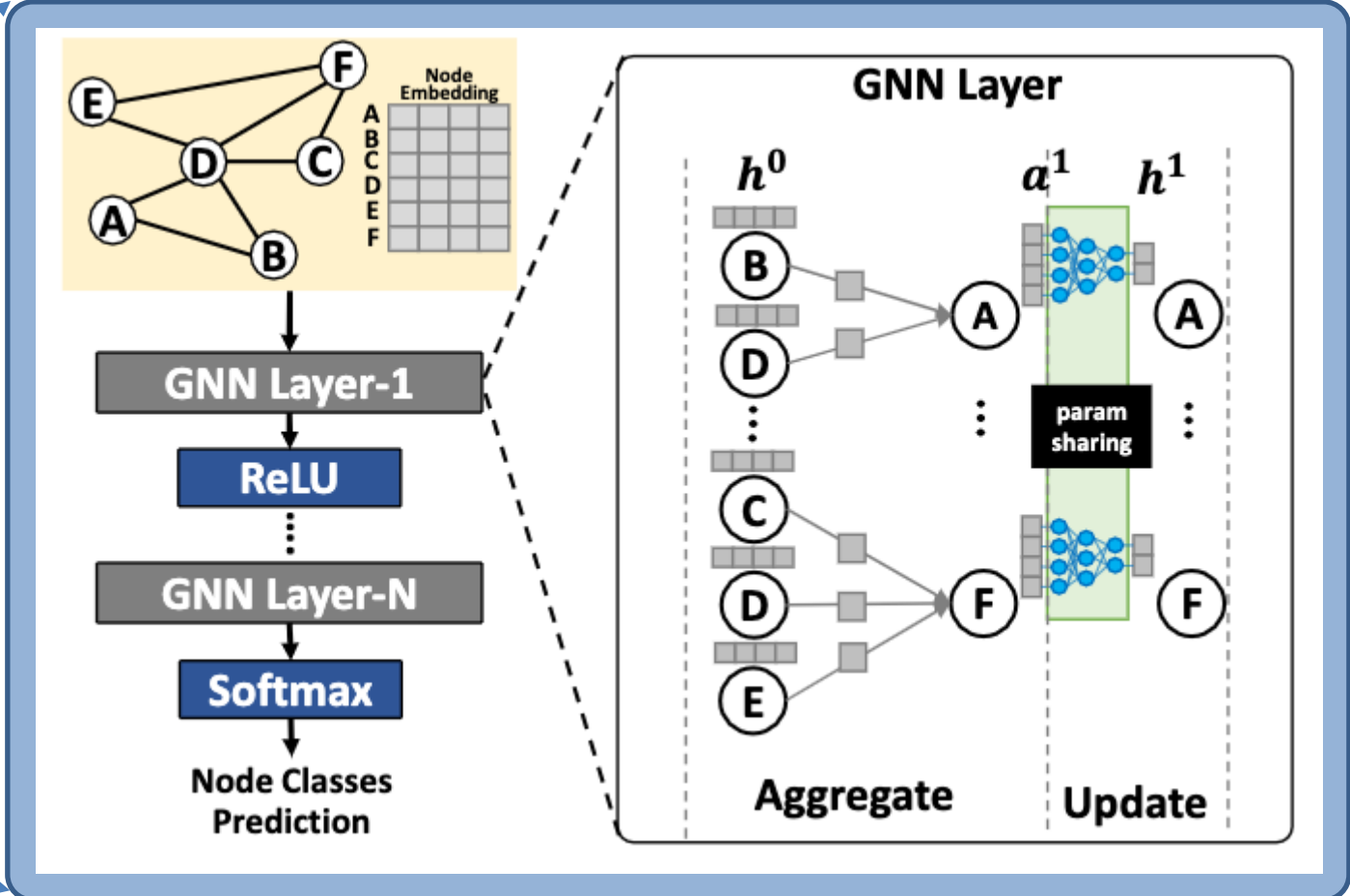
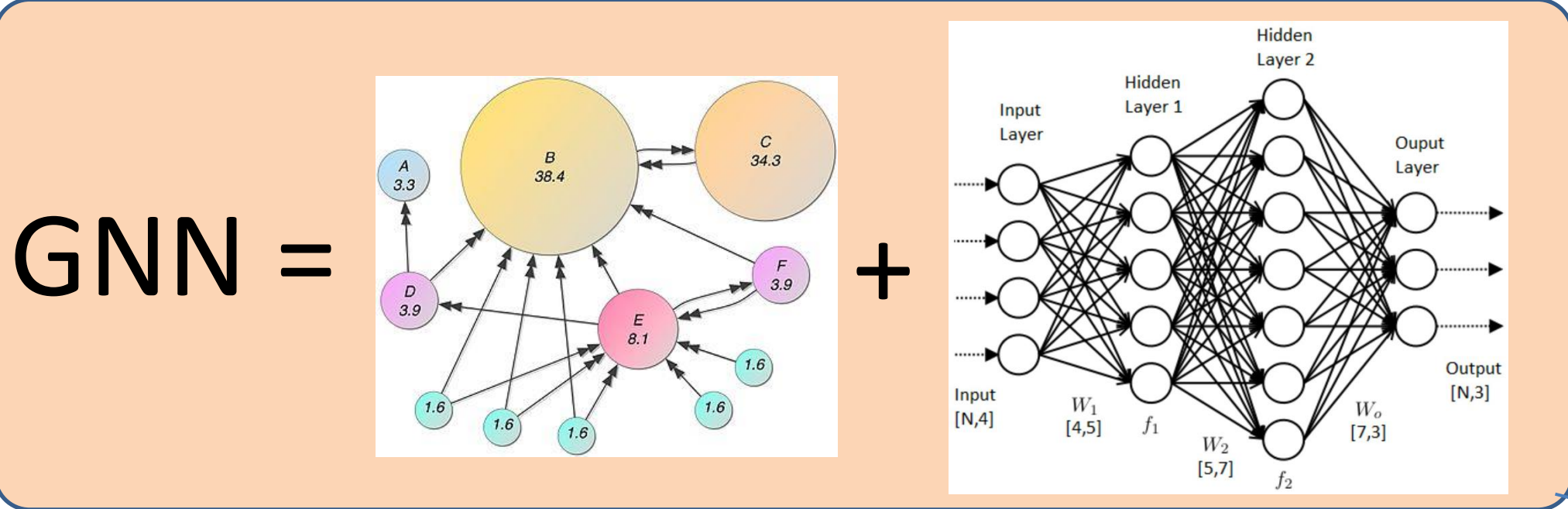


RICE UNIVERSITY

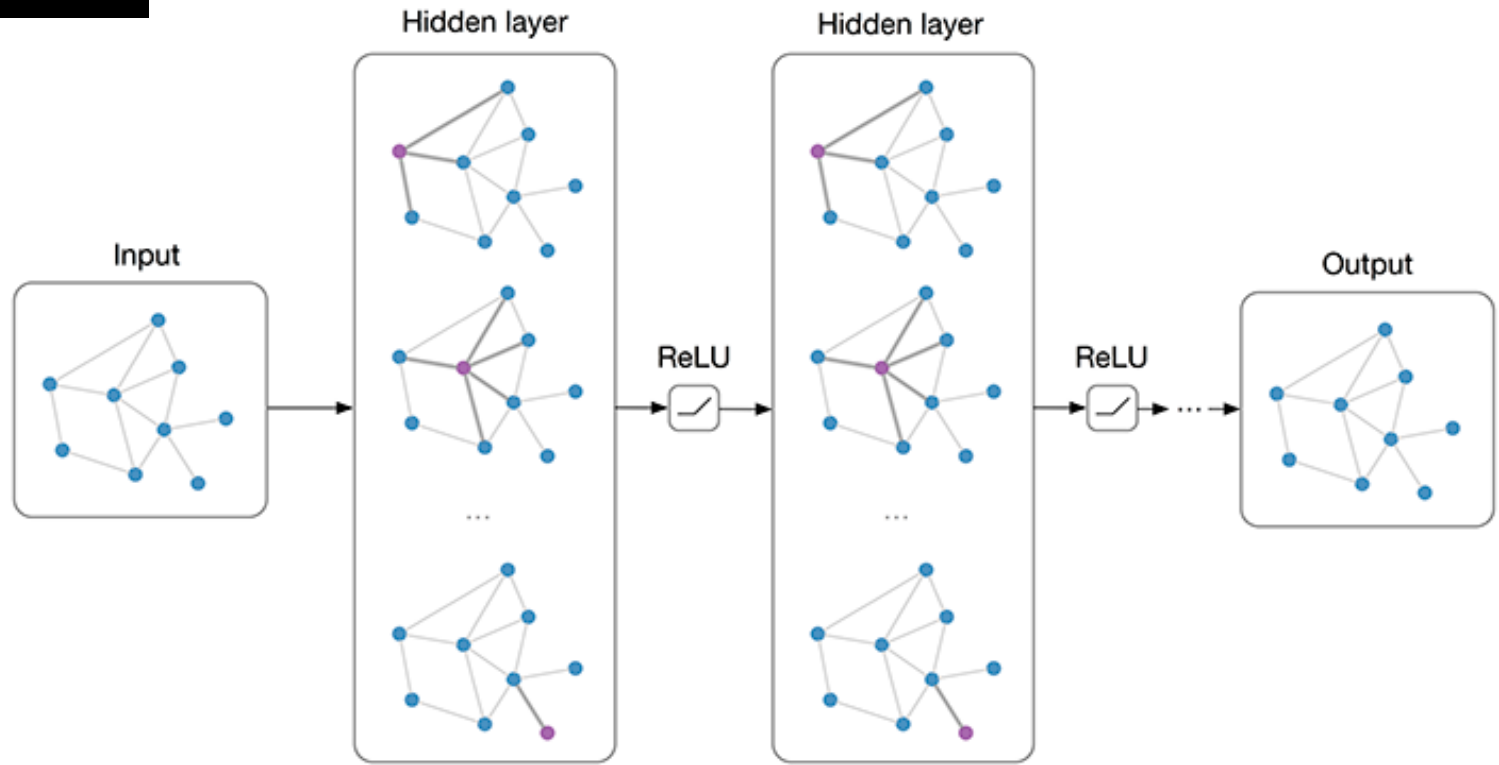
Week-8: Graph Neural Networks

Yuke Wang

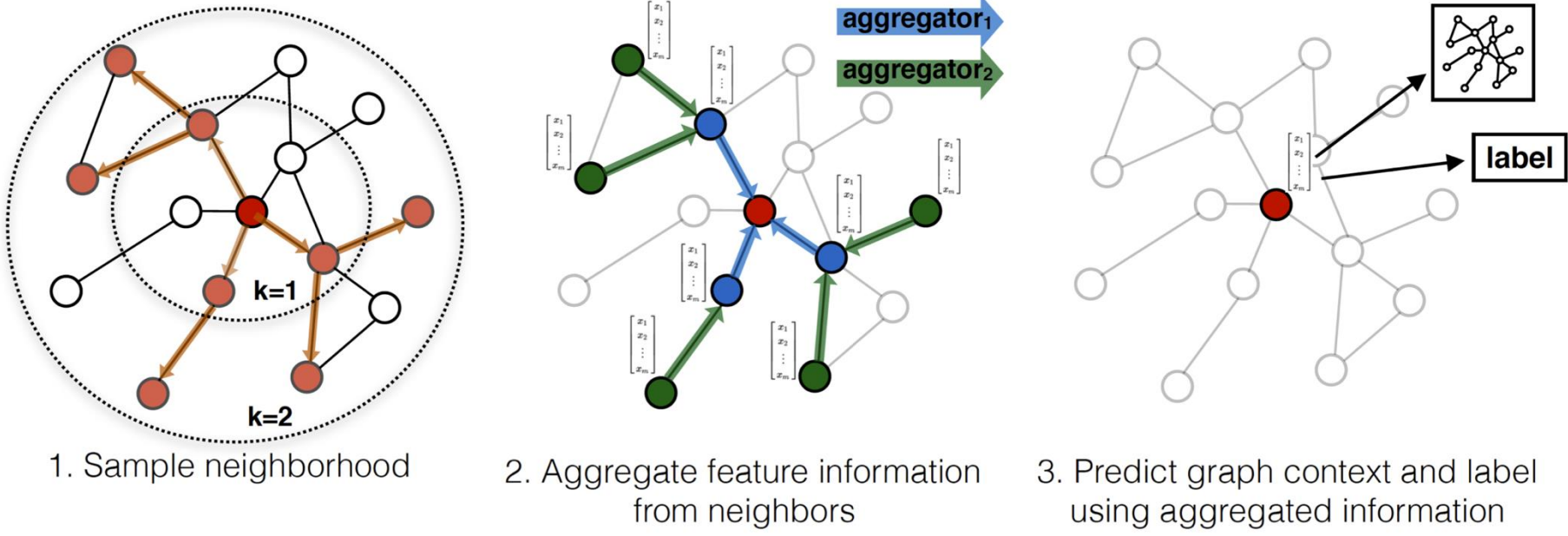
GNN is the key for graph learning



GCN



GraphSAGE



GNN: Challenges

- **Diverse GNN Input:**

- Real-world graphs are huge in size, such as billions/trillions of nodes/edges.
- Real-world graphs have highly irregular structure (random edge connections).

Social Network graphs, such as Twitter graph from SNAP[1], consists of 476 million tweets.

Real-world graphs usually demonstrate highly-scatter distribution of node degrees.

- **Diverse GNN Computation:**

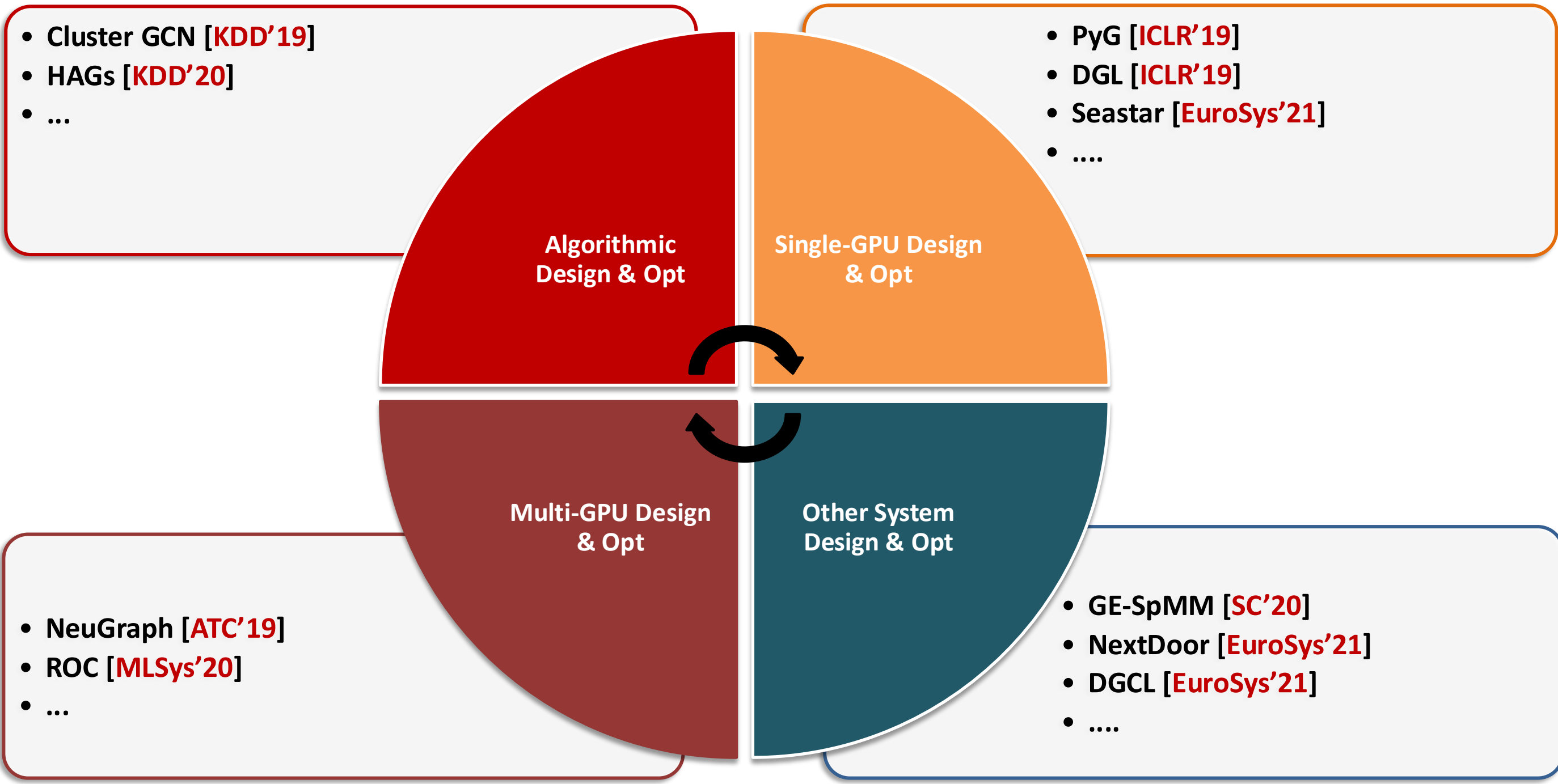
- GNN sparse computation (scatter-and-gather) is intrinsically inefficient for memory access and ill-suited for parallelization on modern platforms.
- GNN computation is an iterative process, with non-trivial computation cost and data movements.

For an graphs with E edges and each nodes with D dimension, the time complexity of a single SAG operation is high as $O(ED)$.

General GNN training processs consists of several hunderd iterations to converge.

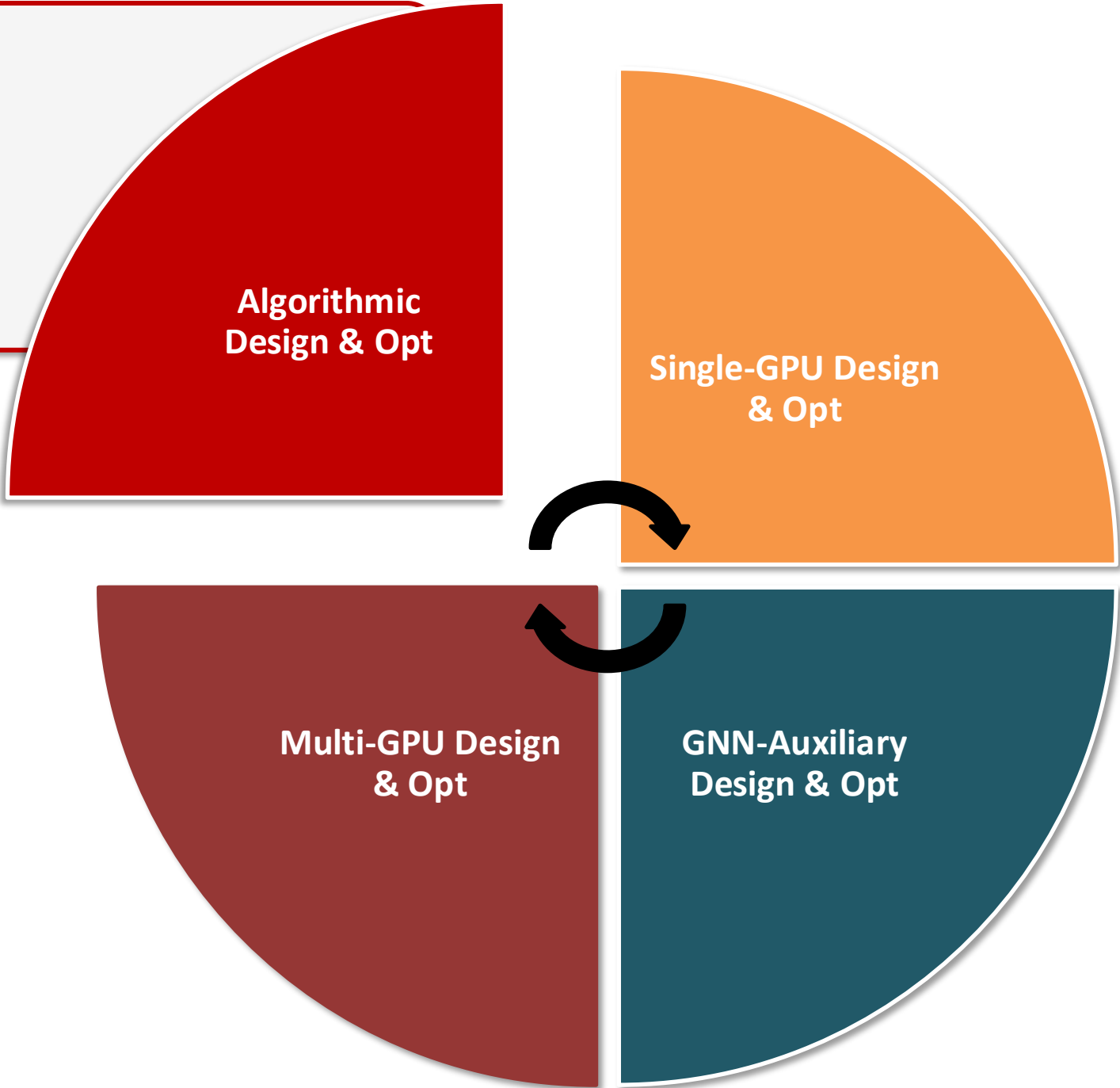
[1] SNAP Datasets: Stanford Large Network Dataset Collection.

Overview



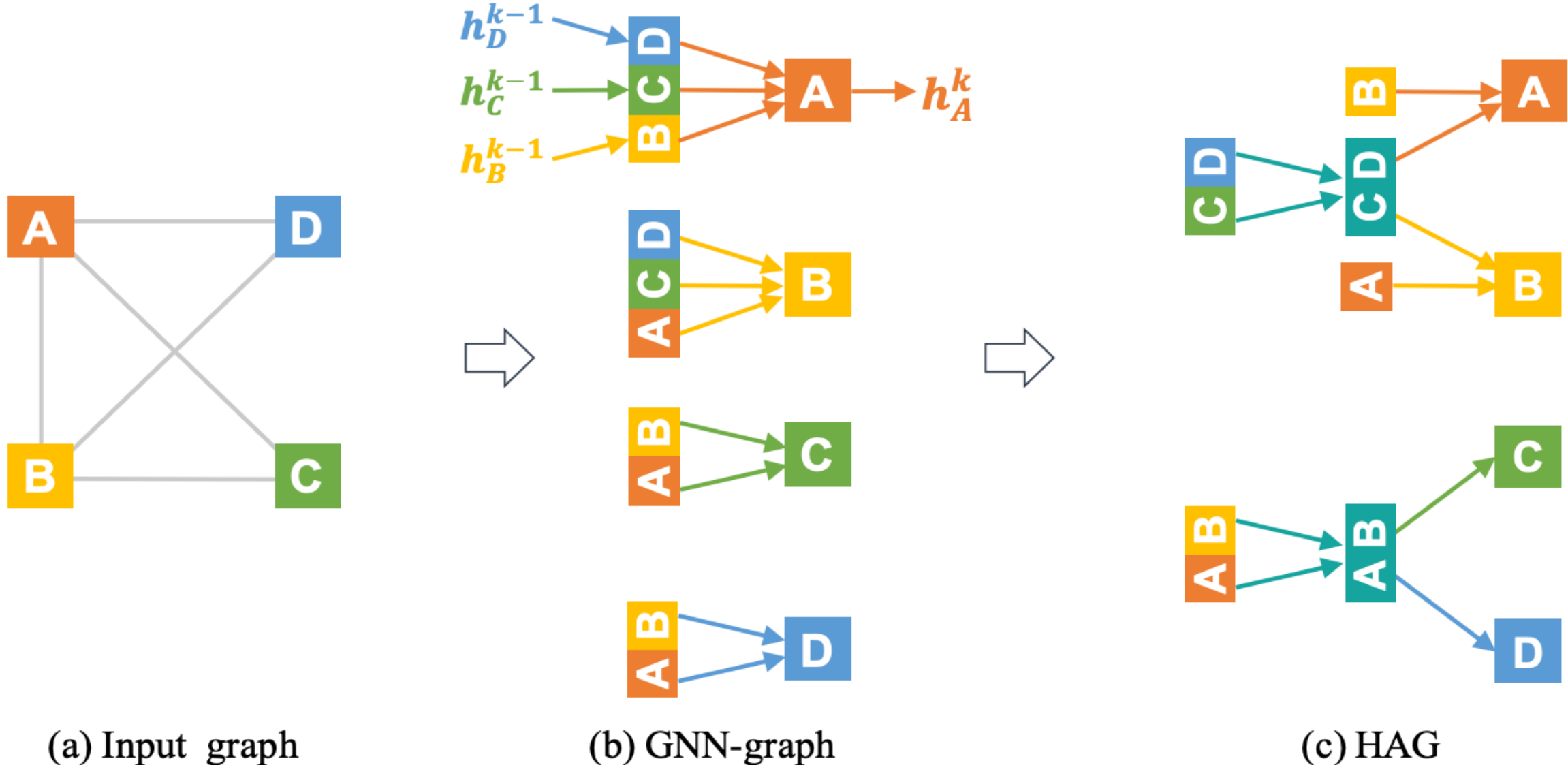
Algorithmic Design & Optimization

- Cluster GCN [KDD'19]
- HAGs [KDD'20]
- ...



HAG

Jia, et al. *Redundancy-free computation graphs for graph neural networks*. KDD'20.



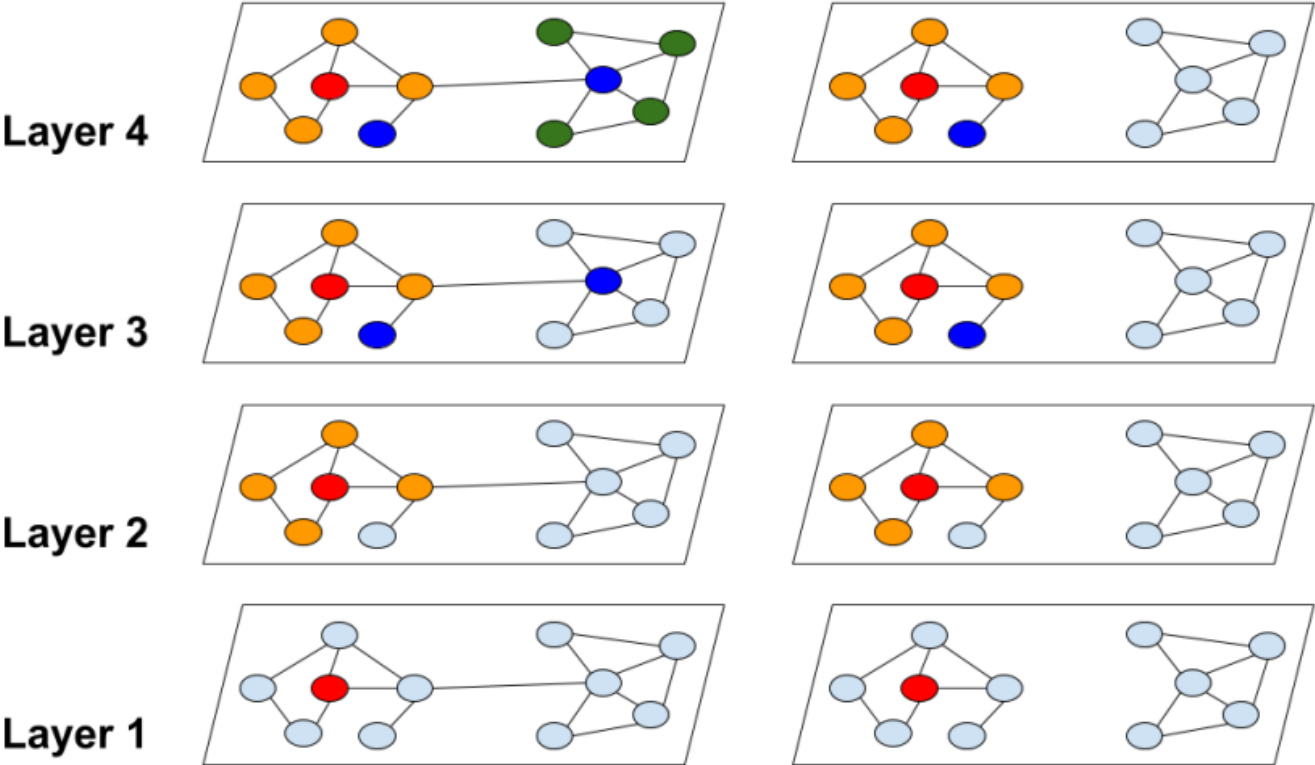
- **Problem:** common neighbors are shared between different nodes, this leads to repeated and inefficient computations.
- **Highlight:** A new GNN graph representation that explicitly avoids redundancy by managing intermediate aggregation results hierarchically.

Increasing the end-to-end training throughput by up to 2.8x and reducing the aggregations and data transfers in GNN training by up to 6.3x and 5.6x, while maintaining the original model accuracy.

Cluster GCN

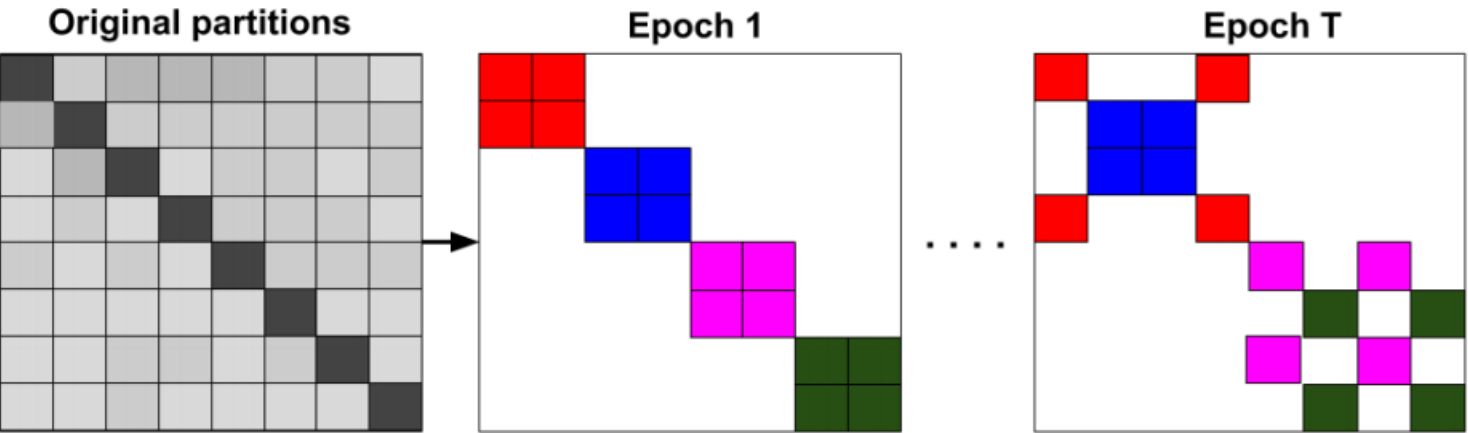
Chiang, et al. *Cluster-gcn: An efficient algorithm for training deep and large graph convolutional networks*. SIGKDD'2019.

➤ Cluster-based Neighbor Expansion.



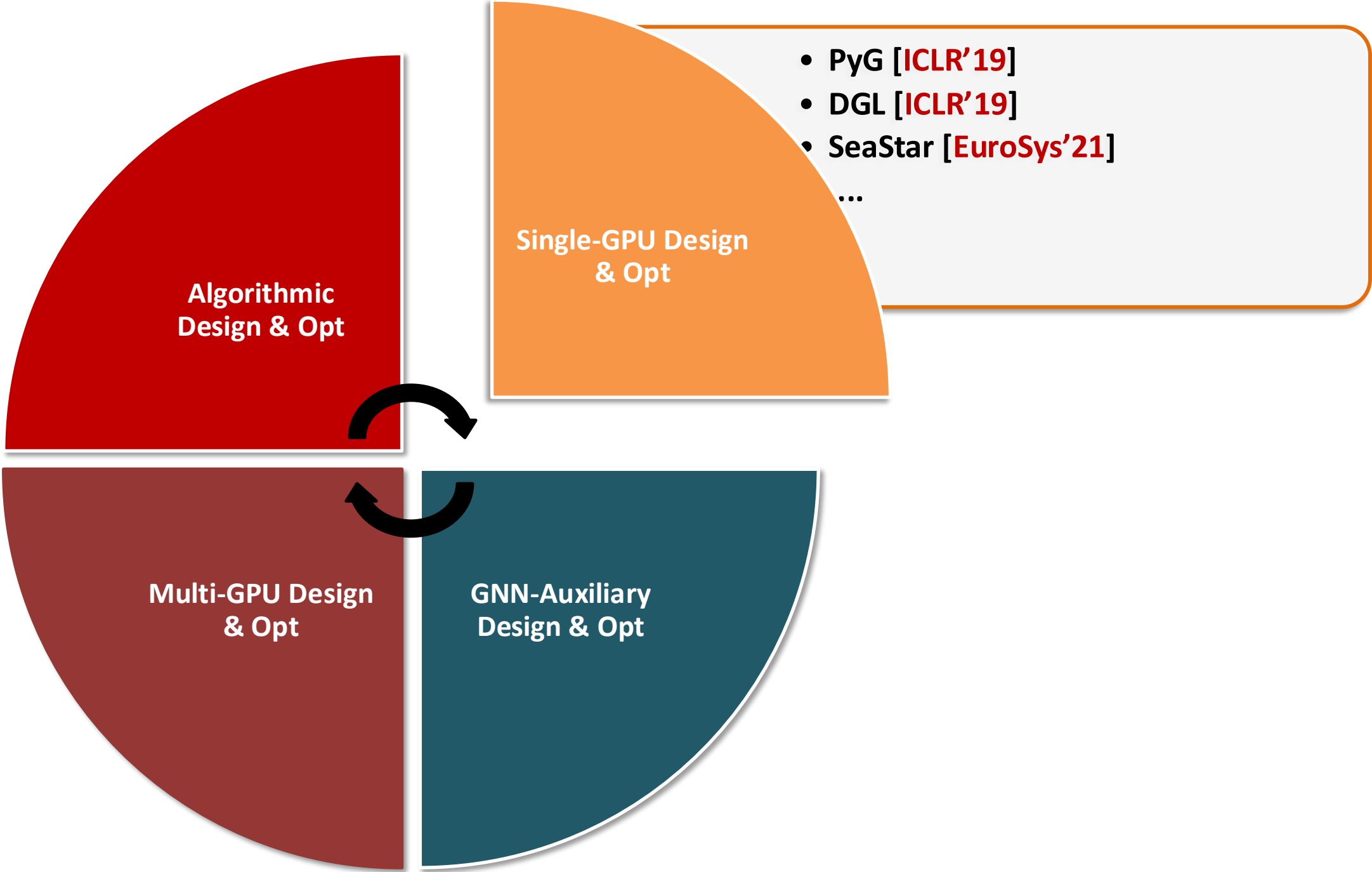
- **Problem:** Existing GNN training algorithms suffer from a high computational cost for keeping the entire graph and the embedding of each node in memory.
- **Highlight:** A novel GCN algorithm that is suitable for SGD-based training by exploiting the graph clustering structure.

➤ Stochastic Multiple Partitions Scheme.

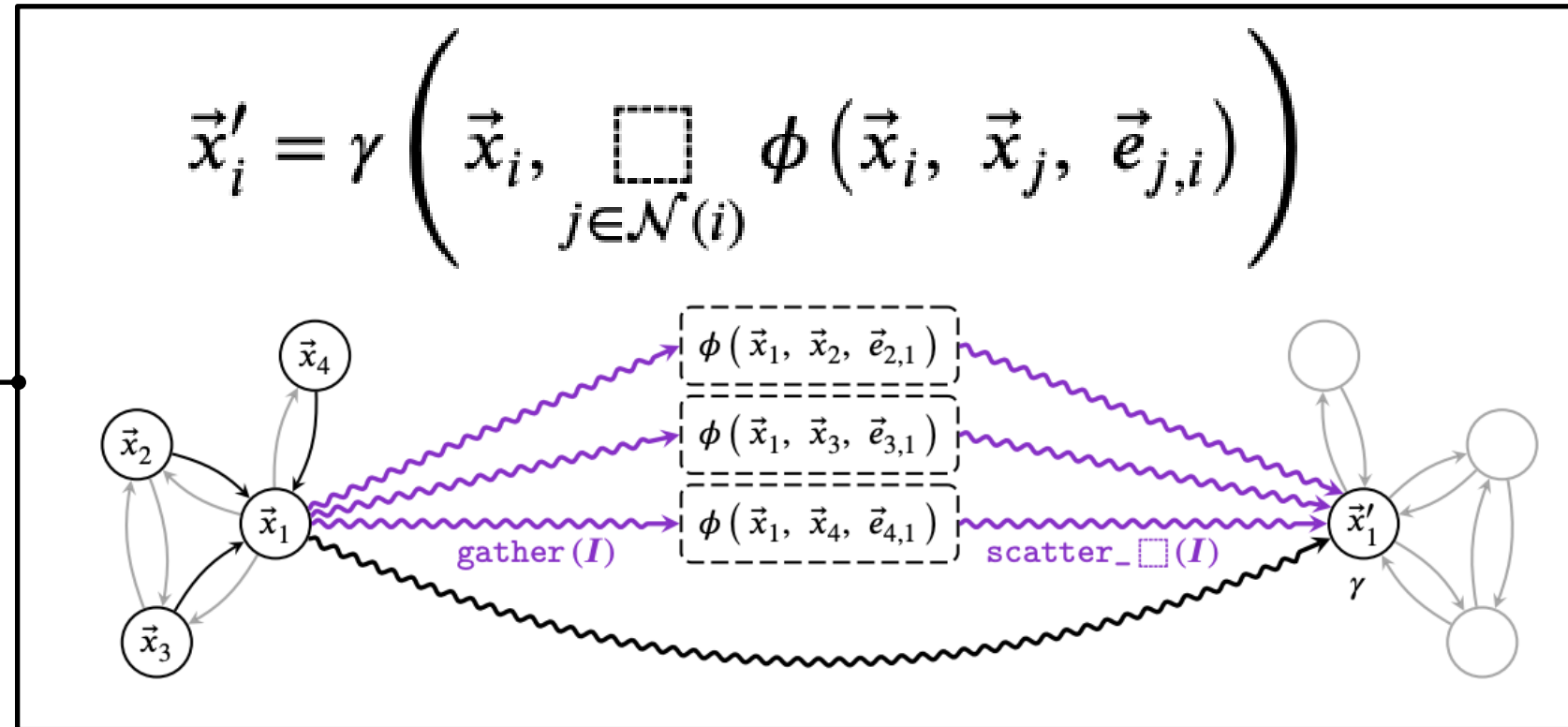


Training a 3-layer GCN on Amazon2M with 2 million nodes and 61 million edges (5x of Reddit), Cluster-GCN is faster than the previous state-of-the-art VR-GCN (1523 seconds vs 1961 seconds) and using much less memory (2.2GB vs 11.2GB).

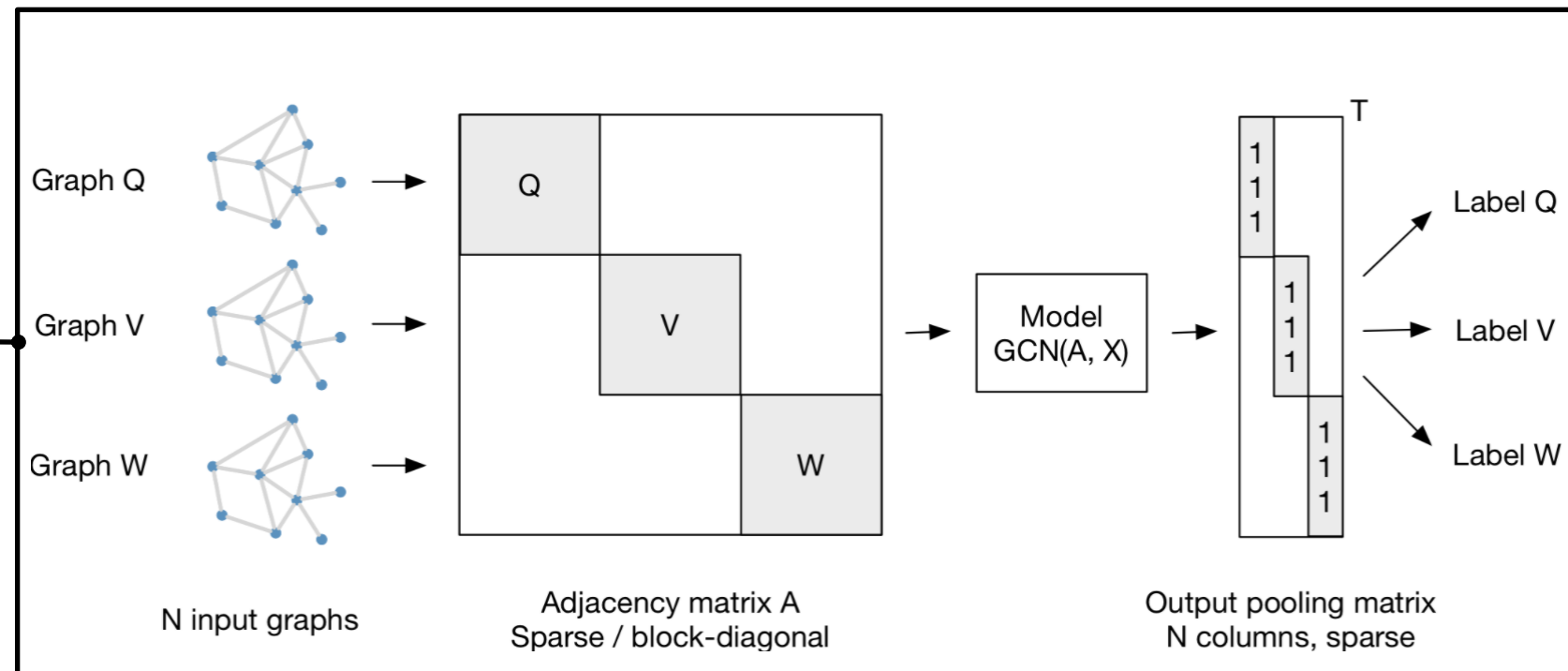
Single-GPU Design & Optimization



1. Programming Model.



2. Mini-Batch Handling.



Deep Graph Library

Wang, et al. *Deep graph library: A graph-centric, highly-performant package for graph neural networks.* ICLR'19.

1. Programming Model.

$$\text{Edge-wise: } \mathbf{m}_k^{(t)} = \phi^e(\mathbf{e}_k^{(t-1)}, \mathbf{v}_{r_k}^{(t-1)}, \mathbf{v}_{s_k}^{(t-1)})$$

$$\text{Node-wise: } \mathbf{v}_i^{(t)} = \phi^v(\mathbf{v}_i^{(t-1)}, \bigoplus_{\substack{k \\ \text{s.t. } r_k=i}} \mathbf{m}_k^{(t)}),$$

2. Graph Storage Management.

Use **DGLGraph** object. Minimize users efforts in choosing graph storage format. (e.g., COO, CSR)

3. Tensorization and Kernel Fusion.

Operate on a batch of edge of features.
Fuse neighbor propagation and SUM reduction.

PyG Vs. DGL

➤ Coverage of GNN models.

		GNet	NGra	Euler	PyG	DGL
Message Passing	arbitrary ϕ^e	✓	✓	✓	✓	✓
	arbitrary ϕ^v	✓	✓	✓	✓	✓
	arbitrary \oplus	✓	✗	✓	✗	✓
Propagation Order	full	✓	✓	✗	✓	✓
	partial	✗	✗	✗	✓	✓
	random walk	✗	✗	✗	✗	✓
	sampling	✗	✗	✓	✗	✓
Graph Type	many & small	✓	✗	✗	✓	✓
	single & giant	✗	✓	✓	✗	✓
	dynamic	✗	✗	✗	✗	✓
System	multi-platform	✗	✗	✗	✗	✓

➤ Performance of Execution.

Dataset		Model	Accuracy	Time		Memory	
$ V $	$ E $			PyG	DGL	PyG	DGL
Cora 3K	11K	GCN	81.31 ± 0.88	0.478	0.666	1.1	1.1
		GAT	83.98 ± 0.52	1.608	1.399	1.2	1.1
CiteSeer 3K	9K	GCN	70.98 ± 0.68	0.490	0.674	1.1	1.1
		GAT	69.96 ± 0.53	1.606	1.399	1.3	1.2
PubMed 20K	889K	GCN	79.00 ± 0.41	0.491	0.690	1.1	1.1
		GAT	77.65 ± 0.32	1.946	1.393	1.6	1.2
Reddit 232K	114M	GCN	93.46 ± 0.06	<i>OOM</i>	28.6	<i>OOM</i>	11.7
Reddit-S 232K	23M	GCN	N/A	29.12	9.44	15.7	3.6

Seastar

Wu et al. *Seastar: vertex-centric programming for graph neural networks*. Eurosys'21.

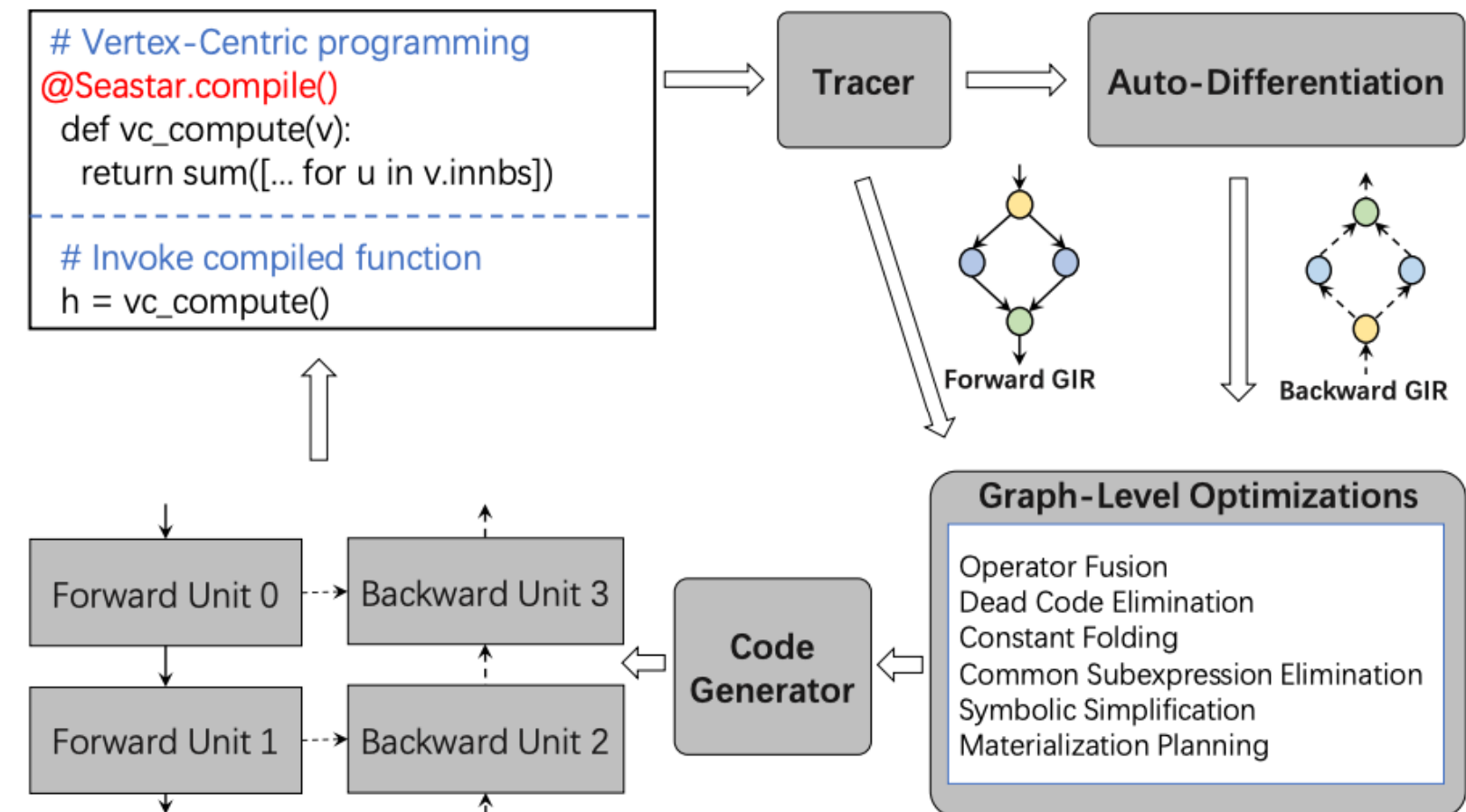
- **Problem:** Existing GNN frameworks use GNN-specific operators plugged into existing deep learning systems, incurring high memory consumption, poor data locality, and large semantic gap between algorithm design and implementation.
- **Highlight:** 1) A vertex-centric programming model for GNN training on GPU. 2) A novel optimizations to produce highly efficient fused GPU kernels for forward and backward passes in GNN training

Seastar achieves up to 2 and 8x less memory consumption, and 14 and 3x faster execution, compared with DGL and PyG.

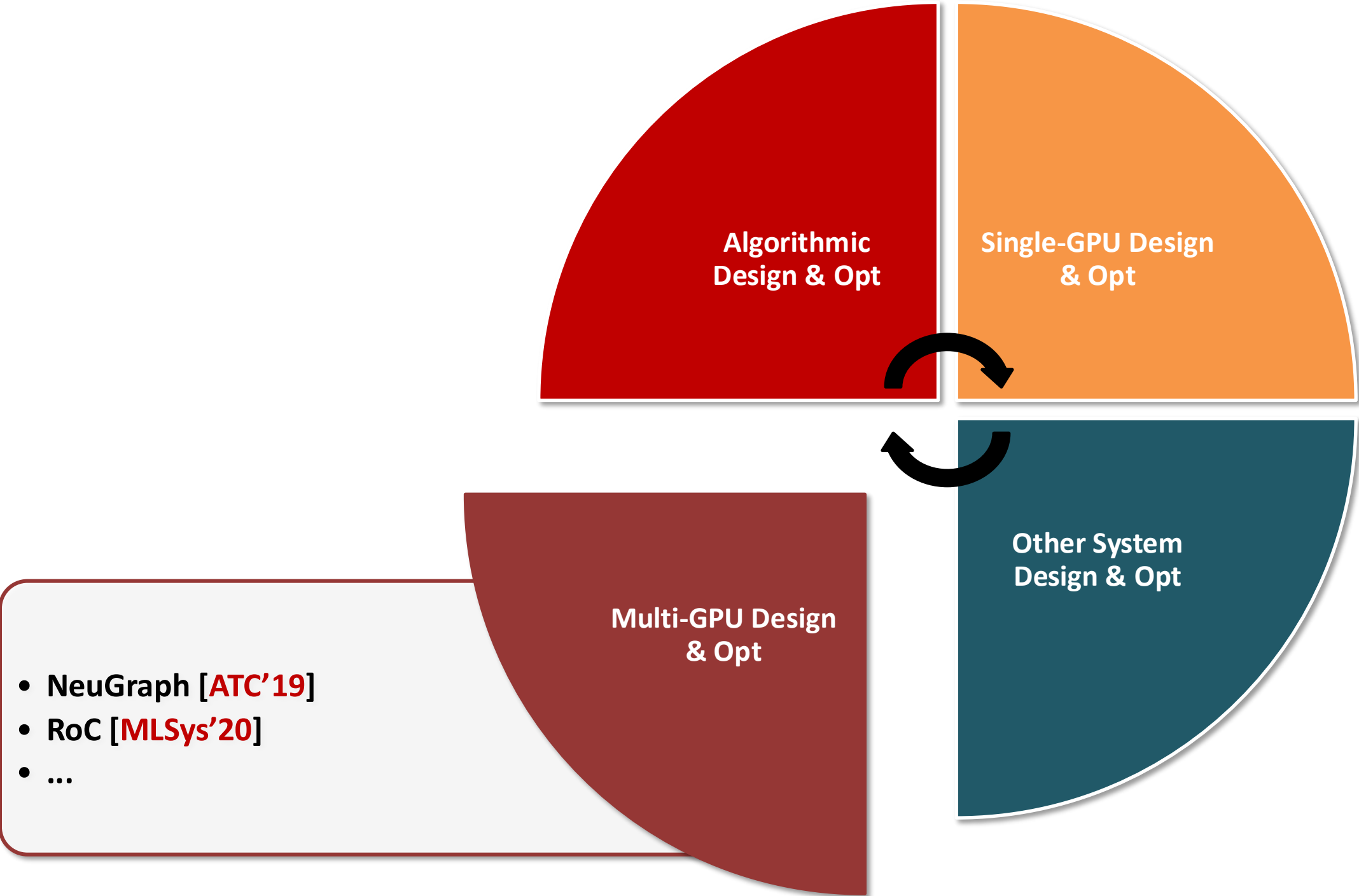
➤ Vertex-Centric Programming.

```
1 #-----Seastar GCN-----#
2 def gcn_forward(self, h, norm):
3     @Seastar.compile(v_feature={'norm':norm, 'h':h})
4     def vc_compute(v):
5         return sum([torch.mm(u.h, self.W) * u.norm
6                     for u in v.innbs])
7     return torch.sigmoid(vc_compute() + self.bias)
```

➤ Overall Optimization and Code Generation Flow.



Multi-GPU Design & Optimization



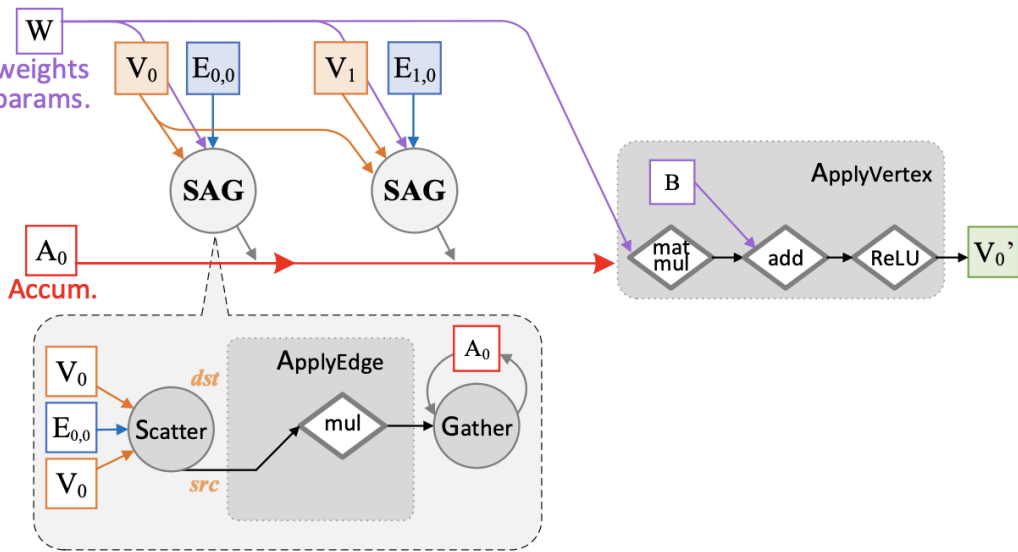
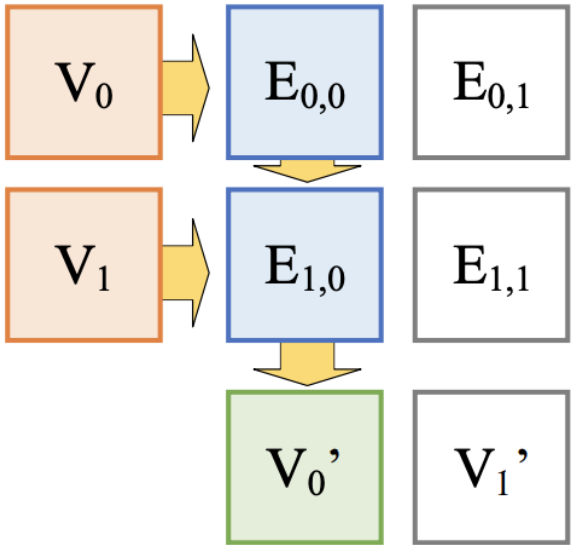
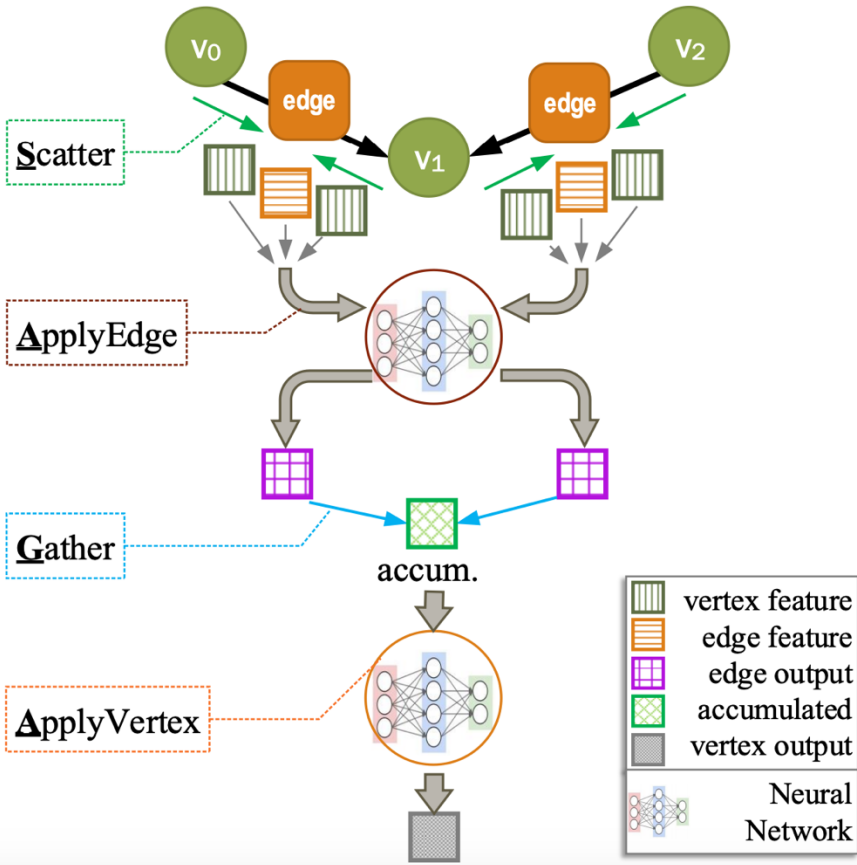
NeuGraph

Lingxiao et al. *NeuGraph: Parallel Deep Neural Network Computation on Large Graphs*. USENIX ATC'19.

- **Problem:** Poor scalability of the existing GNN frameworks on large, densely-connected real-world graphs.
- **Highlight:**
 - A new framework that bridges the graph and dataflow models to support efficient and scalable parallel neural network computation on graphs.
 - Graph computation optimizations into the management of data partitioning, scheduling, and parallelism in dataflow-based deep learning frameworks

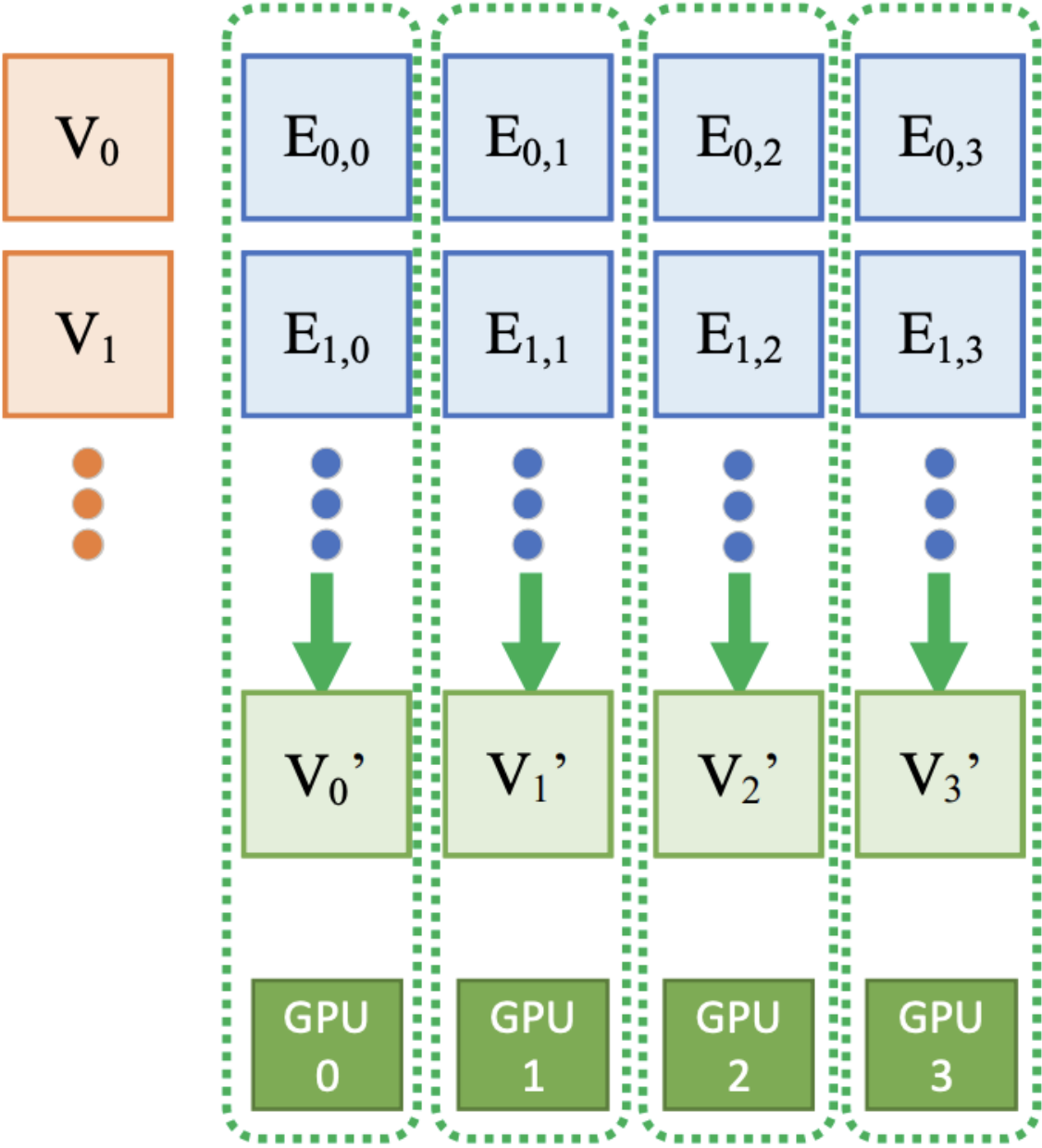
NeuGraph achieves up to 2.5x and 5.0x speedup over the Tensorflow-based GNN design.

SAGA Programming Model & GNN Dataflow Translation.

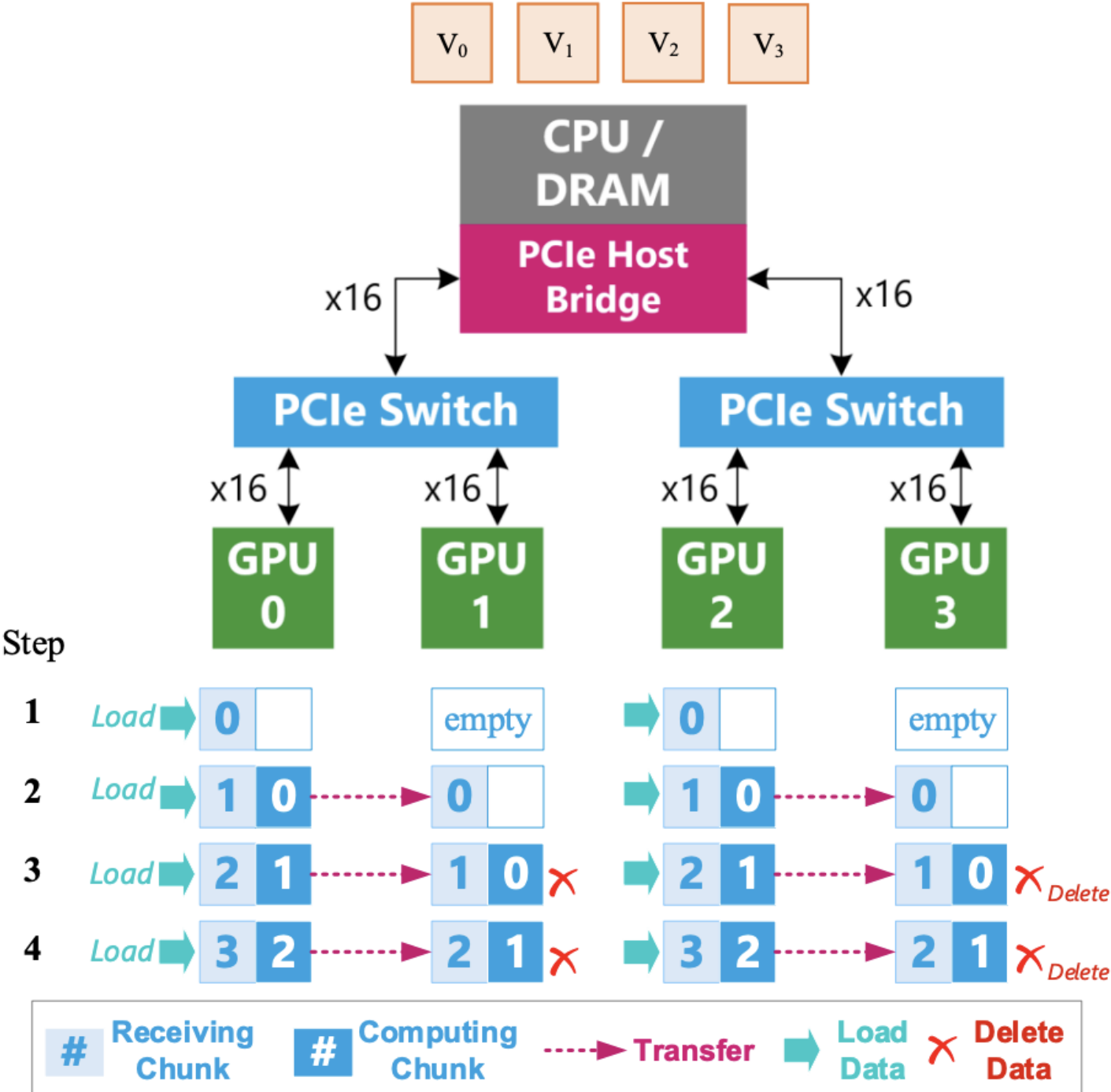


NeuGraph: *Multi-GPU Support*

➤ 2-D Graph Partitioning.



➤ Chain-based Graph Data Forwarding.

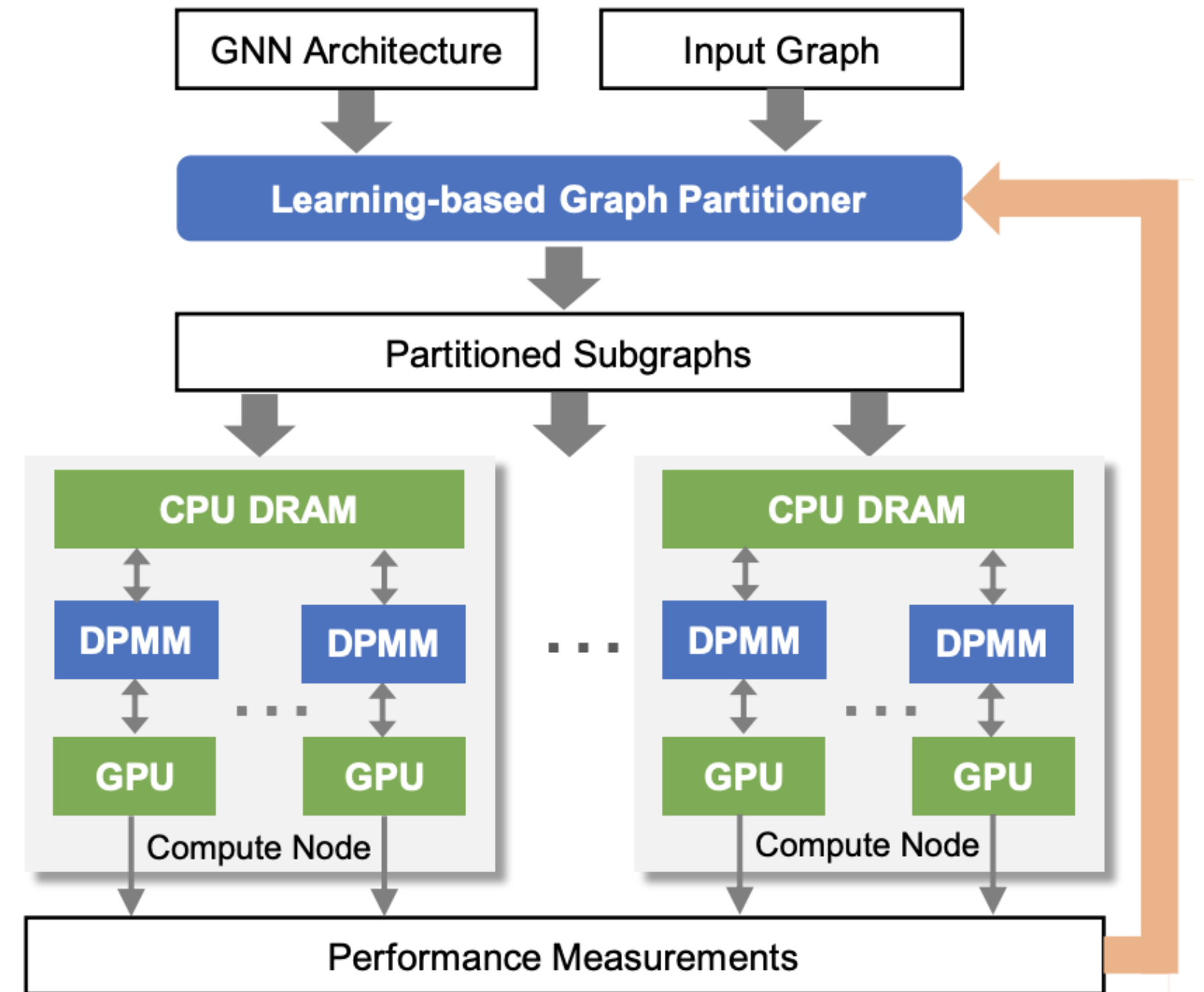


ROC

Z.jia. *Improving the Accuracy, Scalability, and Performance of Graph Neural Networks with Roc*. MLSys'20.

- **Problem:** The current lack of system support has limited the potential application of GNN algorithms on large-scale graphs, and has also prevented the exploration of larger and more sophisticated GNN architectures.
- **Highlight:**
 - A distributed multi-node multi-GPU framework for fast GNN computation.
 - Learning-based graph partitioning.
 - Communication-Minimization Memory Management.
 - Full graph training on larger and deeper GNNs (e.g., 1.5% over sampling method on Reddit).

ROC is up to 4× faster than existing GNN frameworks on a single machine, and can scale to multiple GPUs on multiple machines.



Roc: *Partitioning and Caching*

➤ Cost Model Learning-based Partitioning

We formalize the cost for running a GNN layer l on an input graph \mathcal{G} as follows.

$$t(l, v) = \sum_i w_i(l) x_i(v) \quad (4)$$

$$t(l, \mathcal{G}) = \sum_{v \in \mathcal{G}} t(l, v) = \sum_{v \in \mathcal{G}} \sum_i w_i x_i(v) \quad (5)$$

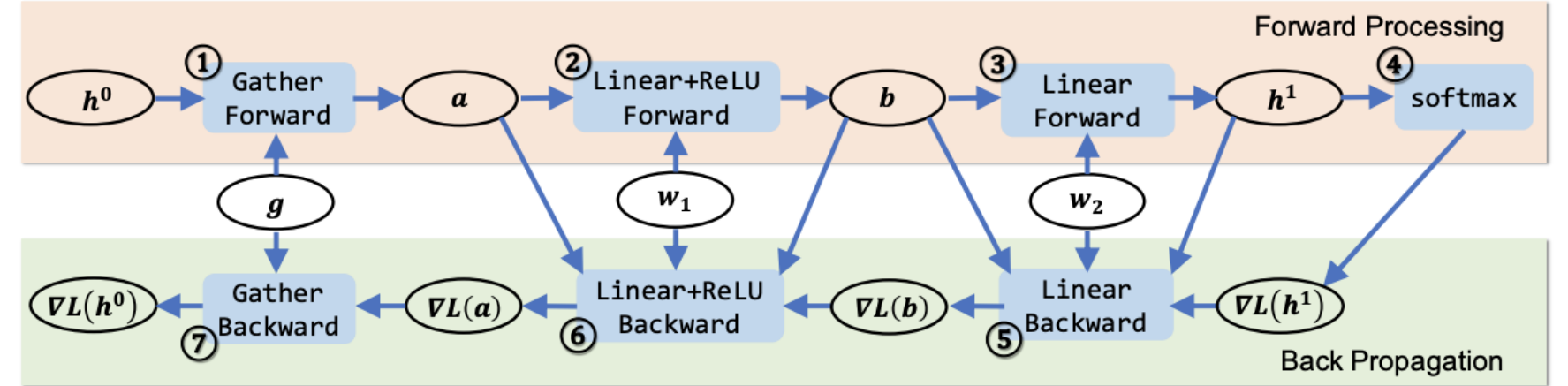
$$= \sum_i w_i \sum_{v \in \mathcal{G}} x_i(v) = \sum_i w_i x_i(\mathcal{G}) \quad (6)$$

where v denotes a vertex in the input graph \mathcal{G} , $w_i(l)$ is a trainable parameter for layer l , $x_i(v)$ is the i -th feature of v , and $x_i(\mathcal{G})$ sums up the i -th feature of all vertices in \mathcal{G} .

Our model minimizes the mean square error over all available data points.

$$Loss(l) = \frac{1}{N} \sum_{i=1}^N (t(l, \mathcal{G}_i) - y(l, \mathcal{G}_i))^2 \quad (7)$$

➤ Cost-Minimization Tensor Caching



Algorithm 1 A recursive dynamic programming algorithm for computing minimum data transfers. $\text{IN}(o_i)$ and $\text{OUT}(o_i)$ return the input and output tensors of the operation o_i , respectively, and $\text{size}(\mathcal{T})$ returns the memory space required to save all tensors in \mathcal{T} .

Table 3. All the valid states and their activation tensors for the GNN architecture in Figure 3.

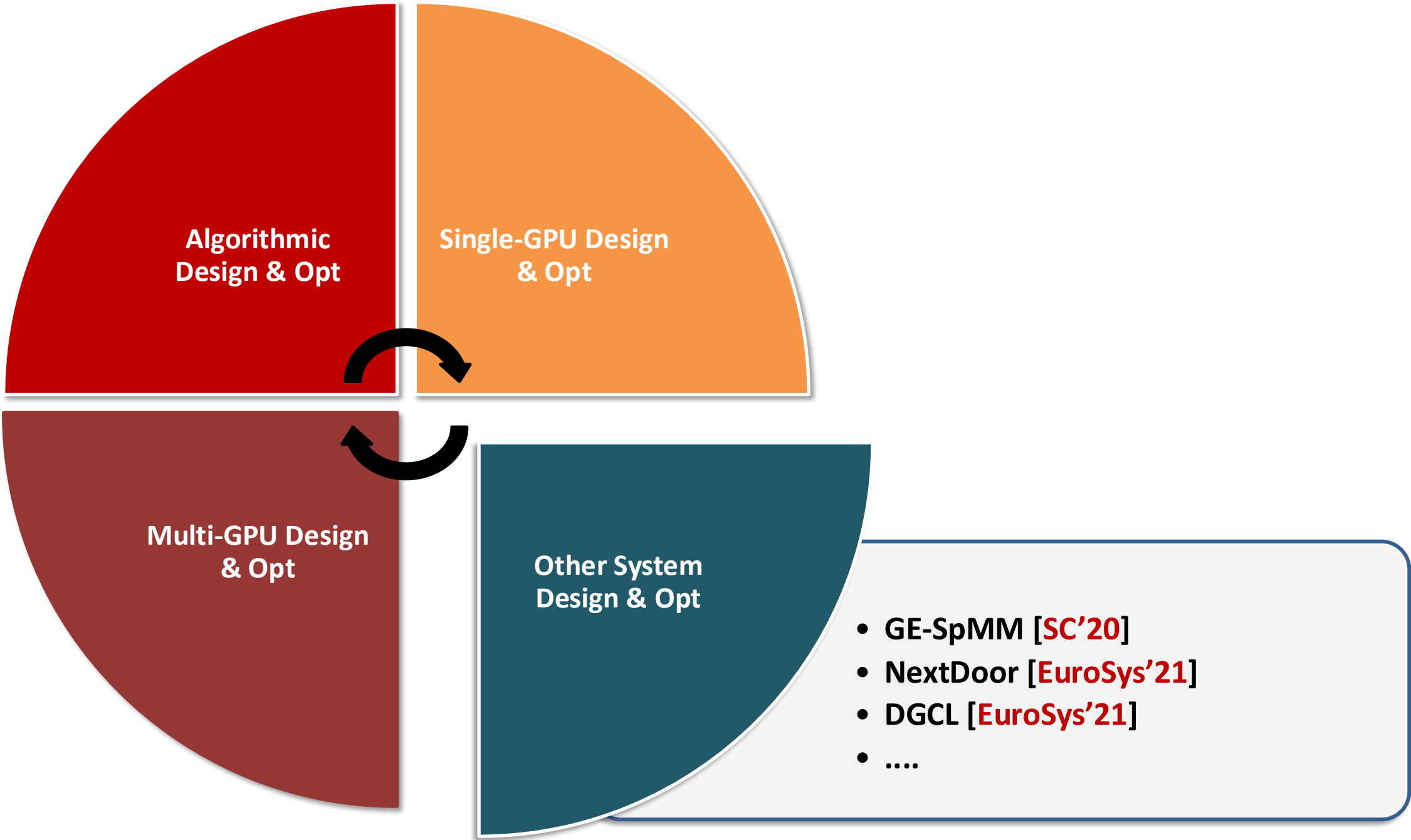
Valid State \mathcal{S}	Activation Tensors $\mathcal{A}(\mathcal{S})$
{①}	{ g, a }
{①, ②}	{ g, a, b, w_1 }
{①, ②, ③}	{ g, a, b, h^1, w_1, w_2 }
{①, ②, ③, ④}	{ $g, a, b, w_1, w_2, \nabla L(h^1)$ }
{①, ②, ③, ④, ⑤}	{ $g, a, b, w_1, \nabla L(b)$ }
{①, ②, ③, ④, ⑤, ⑥}	{ $g, a, \nabla L(a)$ }
{①, ②, ③, ④, ⑤, ⑥, ⑦}	{ }

```

1: Input: An input graph  $g$ , a GNN architecture  $\mathcal{G}$ , and the GPU device memory capacity  $cap$ .
2: Output: Minimum data transfers required to compute  $\mathcal{G}$  on  $g$  within capacity  $cap$ .
3:  $\triangleright \mathcal{D}$  is a database storing all computed COST functions.
4:
5: function COST( $\mathcal{S}, \mathcal{T}$ )
6:   if ( $\mathcal{S}, \mathcal{T}$ )  $\in \mathcal{D}$  then
7:     return  $\mathcal{D}(\mathcal{S}, \mathcal{T})$ 
8:   if  $\mathcal{S}$  is  $\emptyset$  then
9:     return  $\text{size}(\mathcal{T})$ 
10:    $cost \leftarrow \infty$ 
11:   for  $o_i \in \mathcal{S}$  do
12:     if ( $\mathcal{S} \setminus o_i$ ) is a valid state then
13:        $\mathcal{S}' \leftarrow \mathcal{S} \setminus o_i$ 
14:        $\mathcal{T}' \leftarrow (\mathcal{T} \setminus \text{OUT}(o_i)) \cap \mathcal{A}(\mathcal{S}')$ 
15:        $xfer \leftarrow \text{size}(\text{IN}(o_i) \setminus \mathcal{T}')$ 
16:       if  $\text{size}(\mathcal{T} \cup \text{IN}(o_i) \cup \text{OUT}(o_i)) \leq cap$  then
17:          $cost = \min\{cost, \text{COST}(\mathcal{S}', \mathcal{T}') + xfer\}$ 
18:    $\mathcal{D}(\mathcal{S}, \mathcal{T}) \leftarrow cost$ 
19:   return  $\mathcal{D}(\mathcal{S}, \mathcal{T})$ 

```

Other System Design & Optimization



GE-SpMM

Guyue Huang et al. *GE-SpMM: general-purpose sparse matrix-matrix multiplication on GPUs for graph neural networks*. (SC'20)

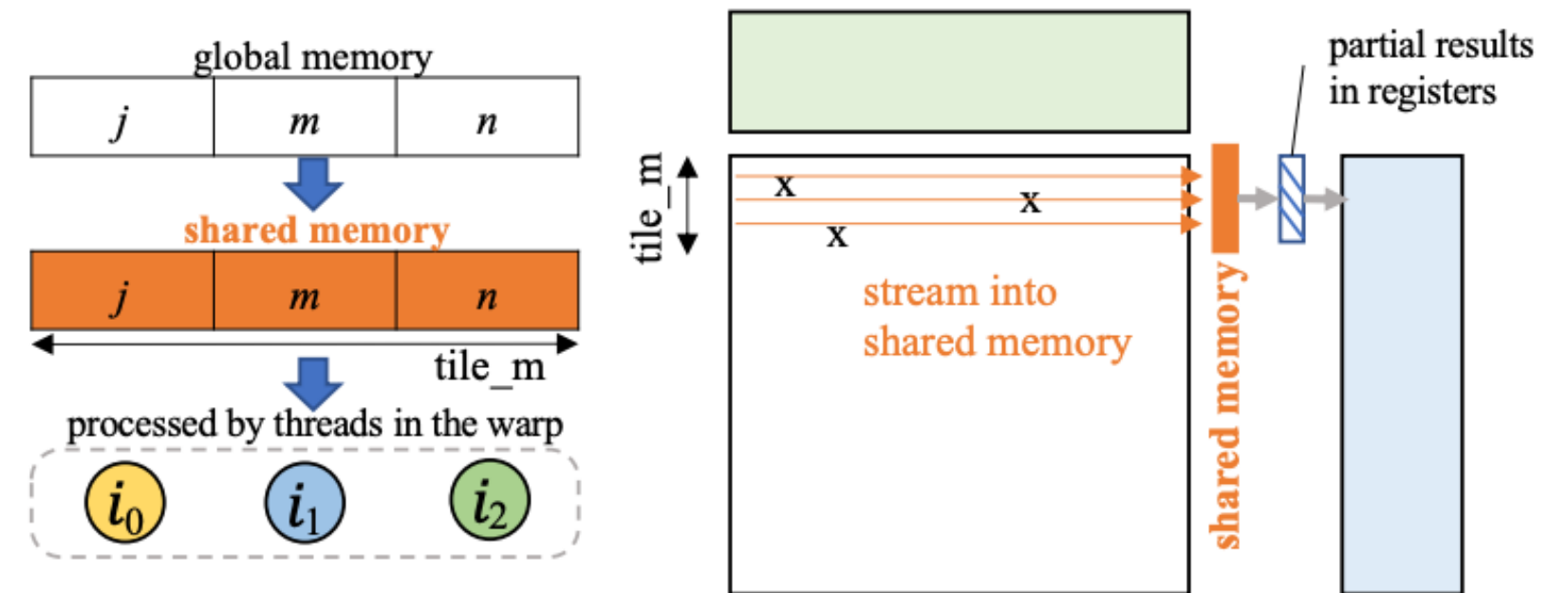
- **Problem:**

- General GNN algorithms require SpMM-like operations (e.g., pooling) are not supported in current high-performance GPU libraries.
- High-overhead data format conversion.
- Lack of SpMM tailored optimizations.

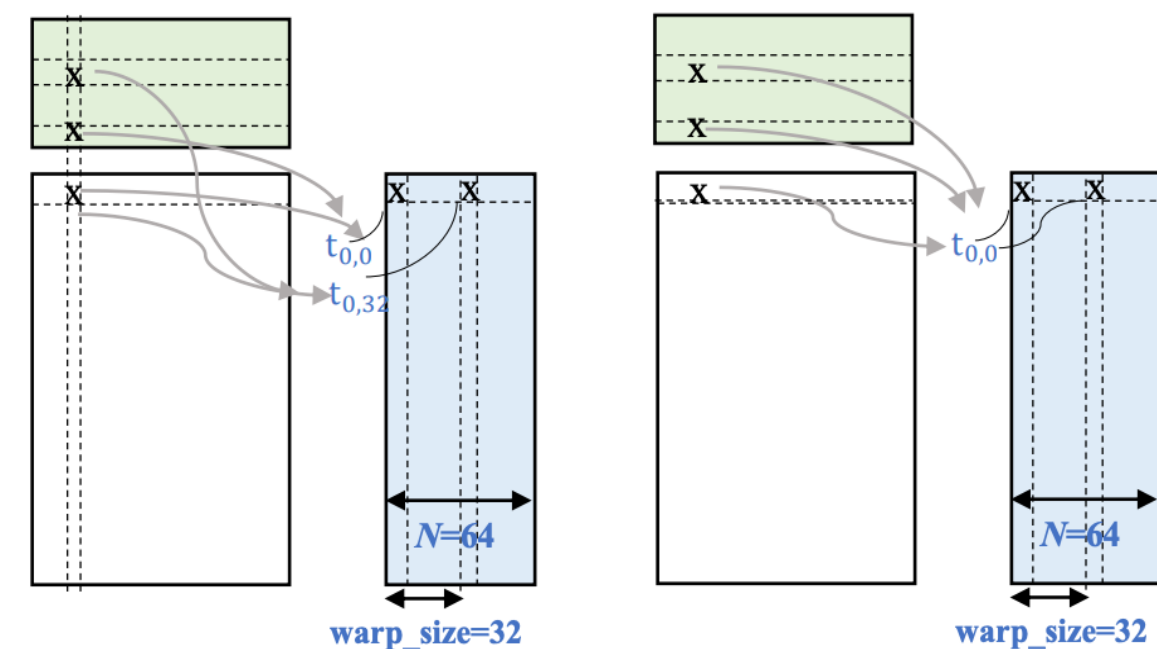
- **Highlight:**

- No preprocessing overheads and support general GNN algorithms.
- Coalesced Row Caching method to process columns in parallel and ensure efficient coalesced access to the GPU global memory.
- Coarsegrained Warp Merging method to reduce redundant data loading among GPU warps.

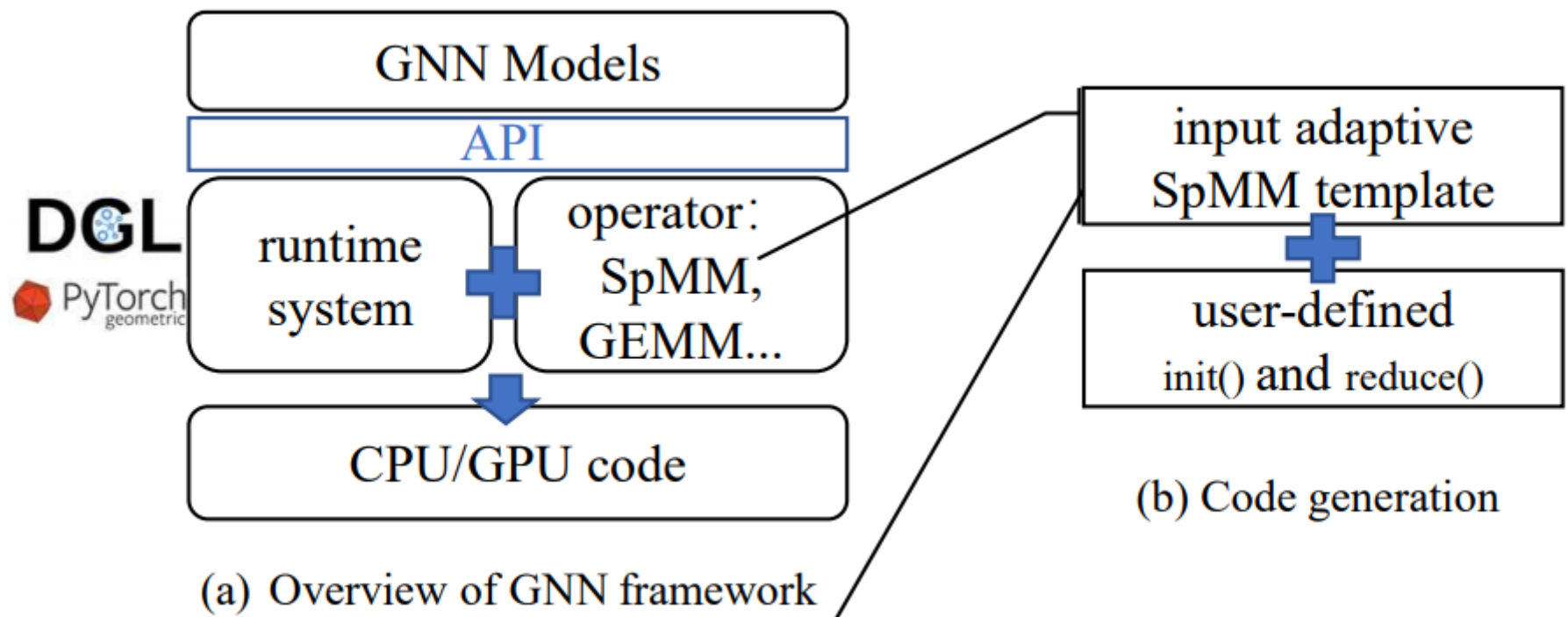
➤ **Coalesced Row Caching**



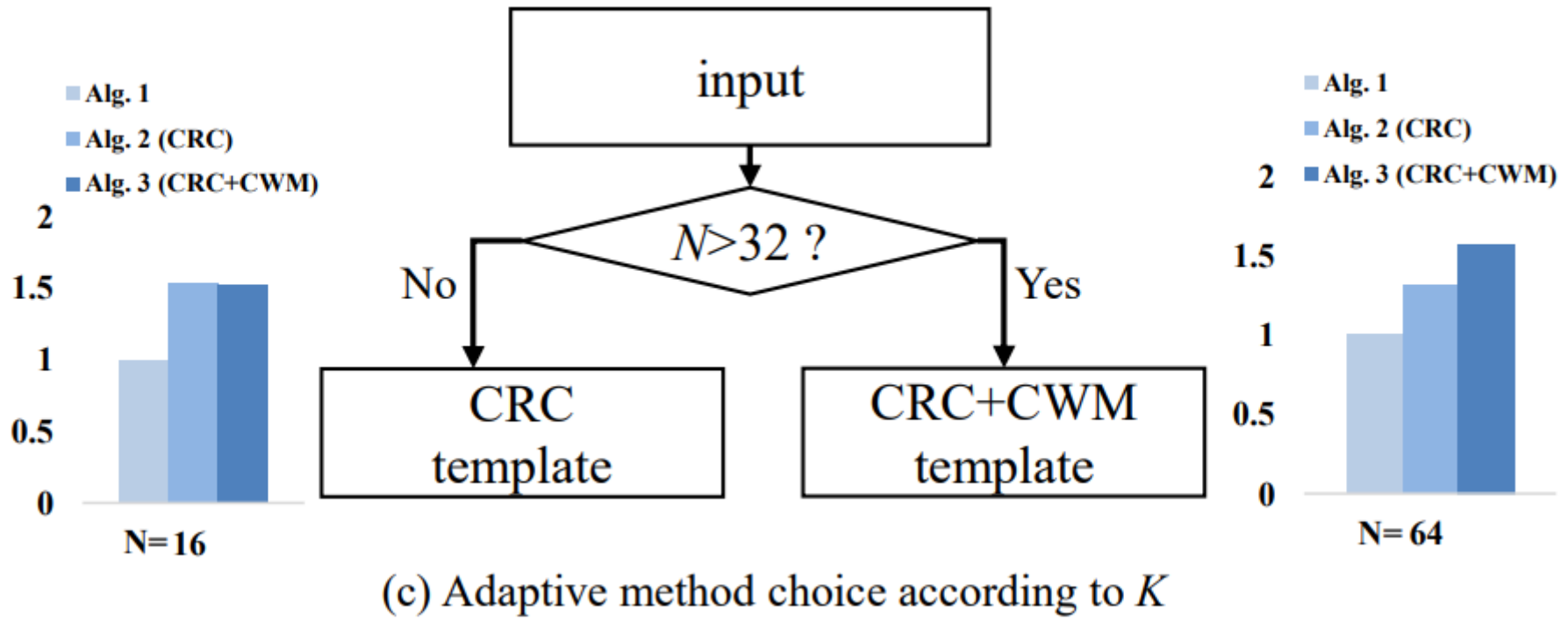
➤ **Coarse-grained Warp Merging**



GE-SpMM: *Overall Flow and Integration*



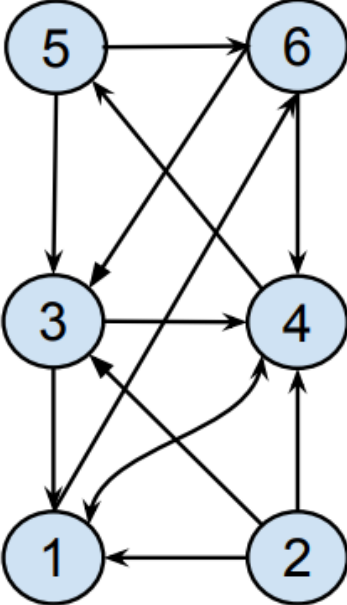
up to $1.41 \times$ speedup over NVIDIA cuSPARSE
 up to $1.81 \times$ over GraphBLAST.
 $3.67 \times$ speedup on GCN and GraphSAGE.



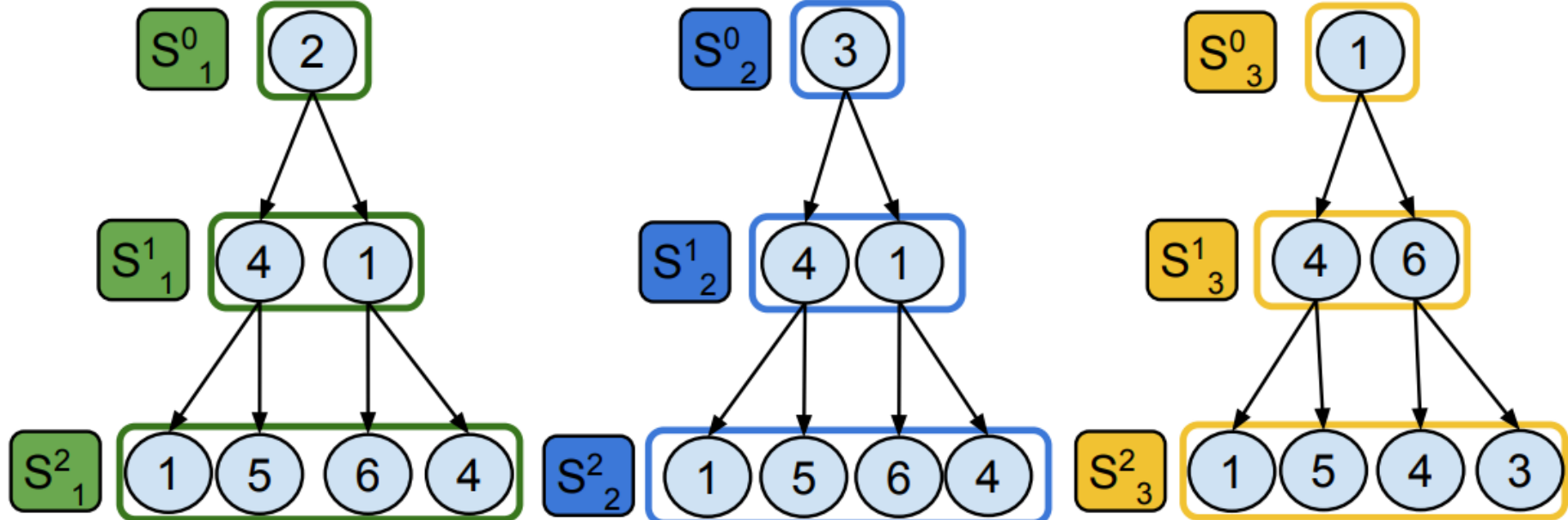
NextDoor

Abhinav Jangda et al. *Accelerating graph sampling for graph machine learning using GPUs*. EuroSys '21

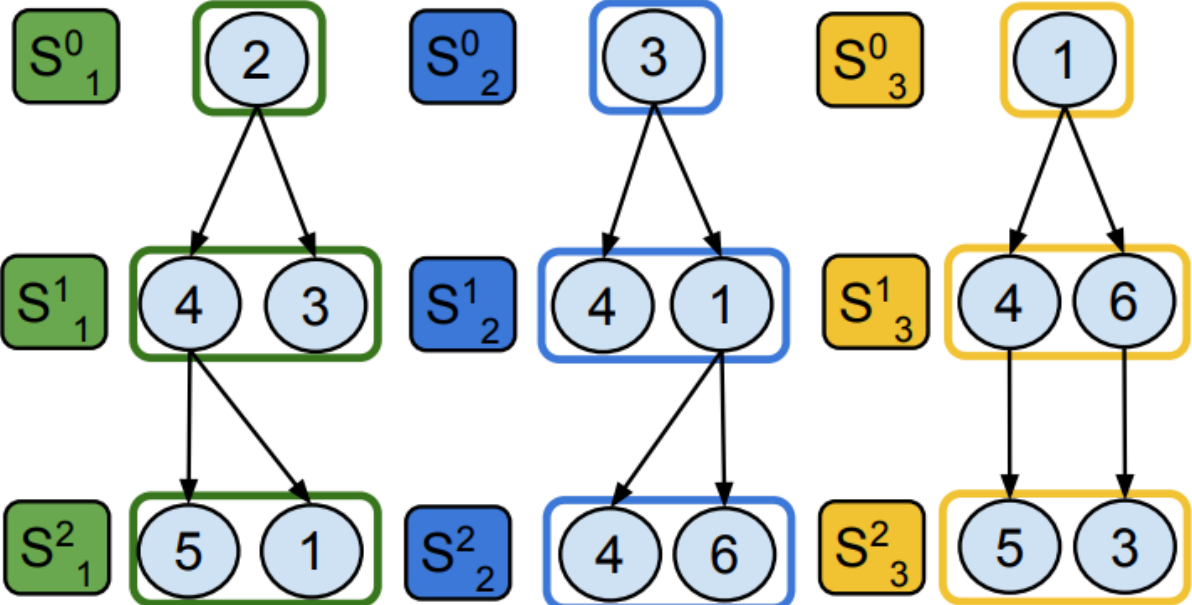
- Problem:** Sampling is an “embarrassingly parallel” problem and may appear to lend itself to GPU acceleration but the irregularity of graphs makes it hard to use GPU resources effectively.
- Highlight:** 1) A system designed to effectively perform graph sampling on GPUs; 2) a new approach to graph sampling that we call transit-parallelism, which allows load balancing and caching of edges; 3) a high-level abstraction for writing a variety of graph sampling algorithms



(a) Graph



(b) 2-hop Neighborhood sampling



(c) Layer Sampling

NextDoor: *Design and Optimizations*

➤ API for GraphSAGE 2-hop sampling.

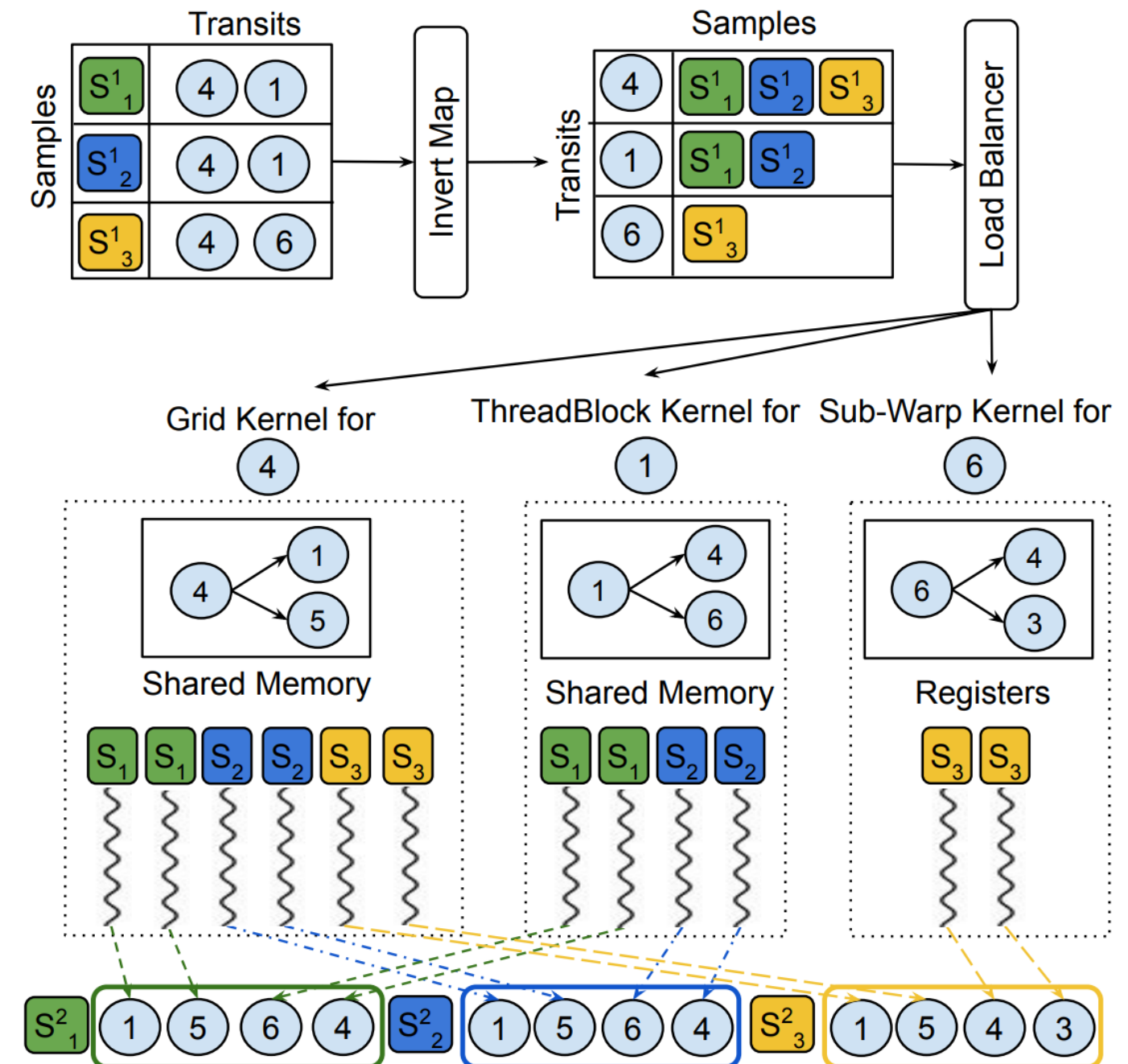
```

1 Vertex next(s, transits, srcEdges, step) {
2   int idx = randInt(0, srcEdges.size());
3   return srcEdges[idx];
4 int steps() {return 2;}
5 int sampleSize(int step) {
6   return (step == 0) ? 25 : 10;}
7 bool unique(int step) {return false;}
8 SamplingType samplingType()
9 {return SamplingType::Individual;}
10 Vertex stepTransits(step, s, transitIdx)
11 {return s.prevVertex(1, transitIdx);}

```

NextDoor provides an order of magnitude speedup over KnightKing (Yang et al, SOSPP'19) for all random walk applications, with speedups ranging from 26.1x to 50x.

➤ Transit Parallel execution of second step of sampling.



DGCL

Cai, Zhenkun, et al. *DGCL: an efficient communication library for distributed GNN training*. Eurosys 2021.

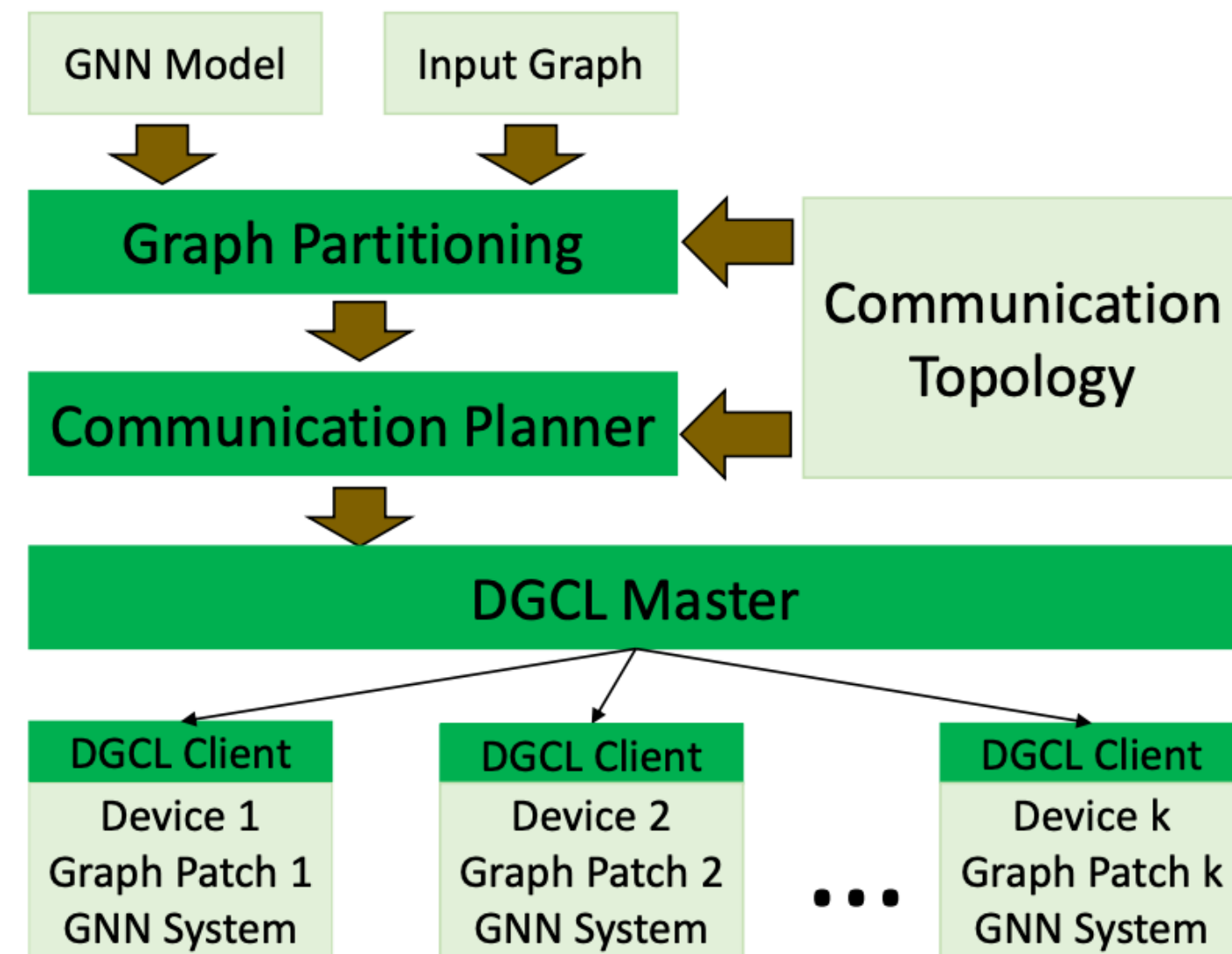
- **Problem:**

- Distributed GNN training, a peer-to-peer communication strategy suffers from high communication overheads.
- Different GPUs require different remote vertex embeddings, which leads to an irregular communication pattern and renders existing communication planning solutions unsuitable.

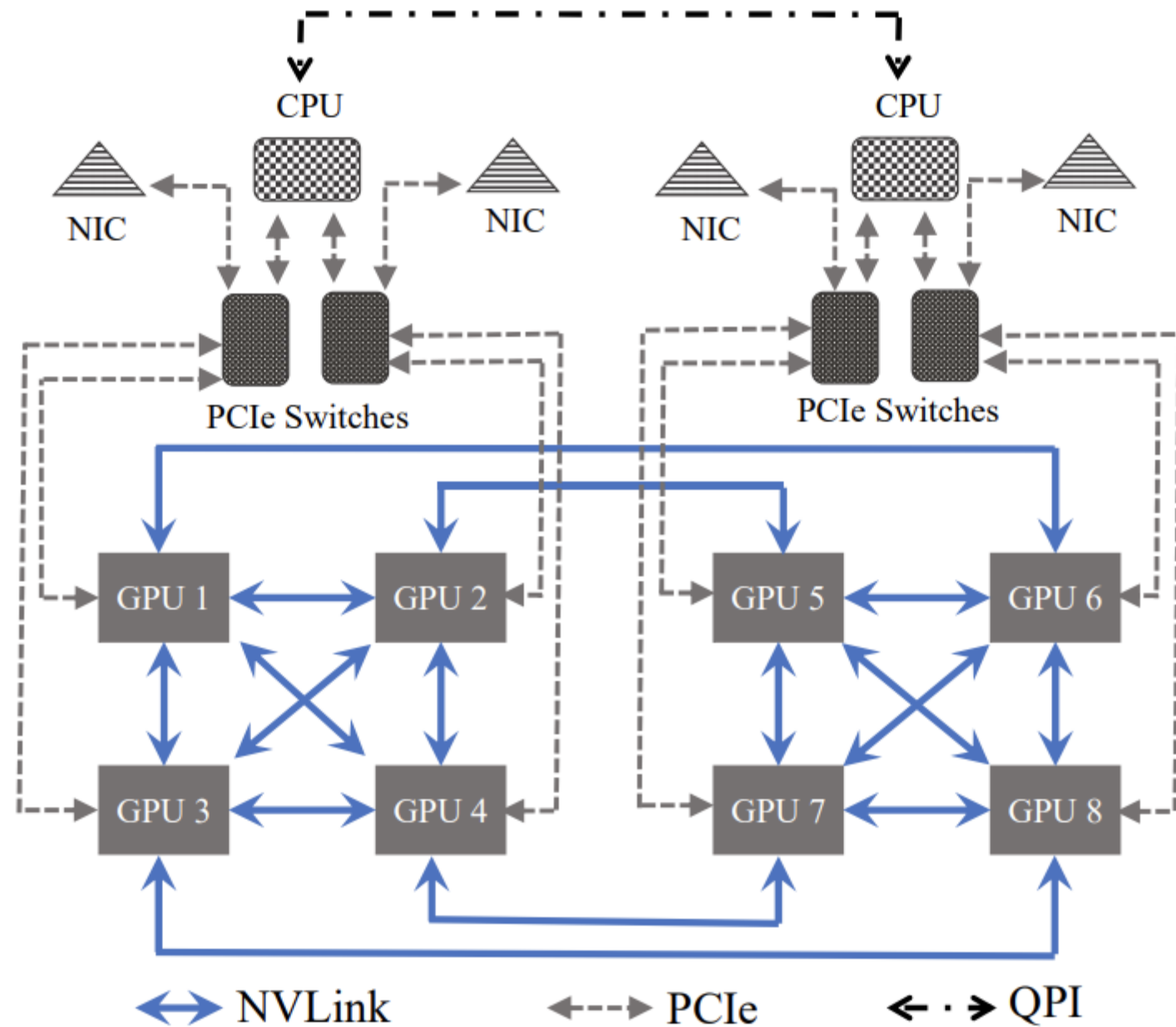
- **Highlight:**

- A distributed graph communication library (DGCL) for efficient GNN training on multiple GPUs.
- A communication planning algorithm tailored for GNN training
- DGCL can be easily adopted to extend existing single-GPU GNN systems to distributed training.

➤ Overview of DGCL



DGCL: *Algorithmic design*

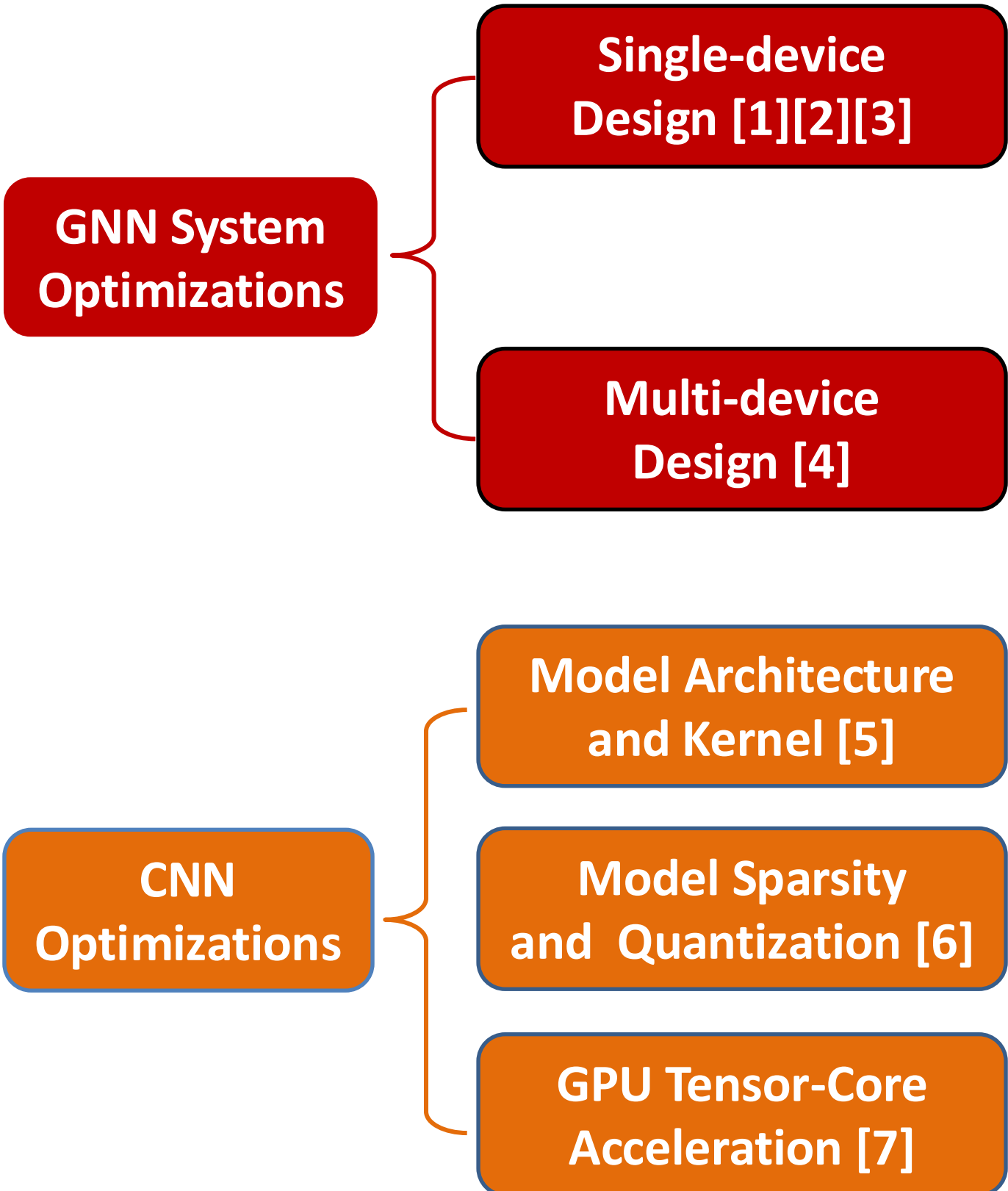


Algorithm 1 Shortest Path Spanning Tree (SPST)

- 1: **Input:** The communication topology graph $D(V', E')$, the source GPU s_u and the destination GPUs D_u for each vertex $u \in V$.
 - 2: **Output:** The communication strategies for all vertices in V .
 - 3: **for** $u \in V$ **do**
 - 4: $N_{src} = \{s_u\}$
 - 5: **for** i from 1 to $|D_u|$ **do**
 - 6: $C = \text{incrmentalLinkCost}(D, S) \triangleright \text{Use Algorithm 2}$
 - 7: Run $\text{dijkstra}(N_{src}, D, C, D_u \setminus N_{src})$
 - 8: Let p be the shortest path among the set of shortest paths from each $d \in N_{src}$ to each $d' \in D_u \setminus N_{src}$ computed in Dijkstra algorithm
 - 9: $N_{src} = N_{src} \cup \{\text{vertices on } p\}$
 - 10: $S = S \cup \{\text{edges on } p\}$
 - 11: **end for**
 - 12: **end for**
 - 13: **Return:** S
-

Reduce the communication time of the peer-to-peer communication by 77.5% on average and the training time for an epoch by up to 47%.

My Research



[1] Y. Wang et al. ***GNNAdvisor: An Efficient Runtime System for GNN Acceleration on GPUs.*** **OSDI'21**

[2] Y. Wang et al. ***TC-GNN: Accelerating Sparse Graph Neural Network Computation Via Dense Tensor Core on GPUs*** (In Submssion)

[3] Y. Wang et al. ***QGTC: Accelerating Quantized GNN via GPU Tensor Core.*** (In Submssion)

[4] Y. Wang et al. ***Accelerating Graph Neural Networks on Multi-GPU Platforms*** (In Progress)

[5] Y. Wang et al. ***DSXplore: Optimizing Convolutional Neural Networks via Sliding-Channel Convolutions.*** **IPDPS'21.**

[6] L. Liu et al. ***Boosting Deep Neural Network Efficiency with Dual-Module Inference*** **ICML'20.**

[7] B. Feng et al. ***APNN-TC: Accelerating Arbitrary Precision Neural Networks on Ampere GPU Tensor Core*** (In Submission)



RICE UNIVERSITY

yuke.wang@rice.edu