

USENIX ATC '19, Renton, WA, USA

NeuGraph:

Parallel Deep Neural Network Computation on Large Graphs

Lingxiao Ma[†], Zhi Yang[†], Youshan Miao[‡], Jilong Xue[‡], Ming Wu[‡], Lidong Zhou[‡], Yafei Dai[†]

[†] *Peking University*

[‡] *Microsoft Research*



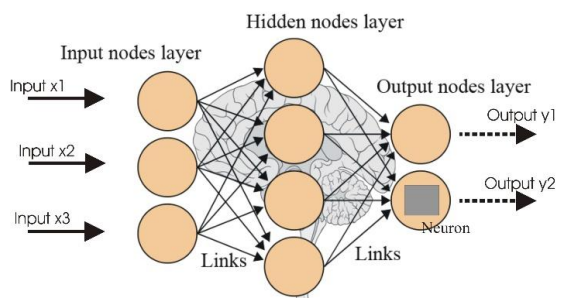


Self-Driving



Personal Assistant

Input Feature Vector



Neural Networks

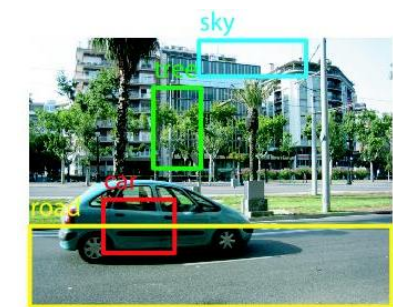


Image Object Detection



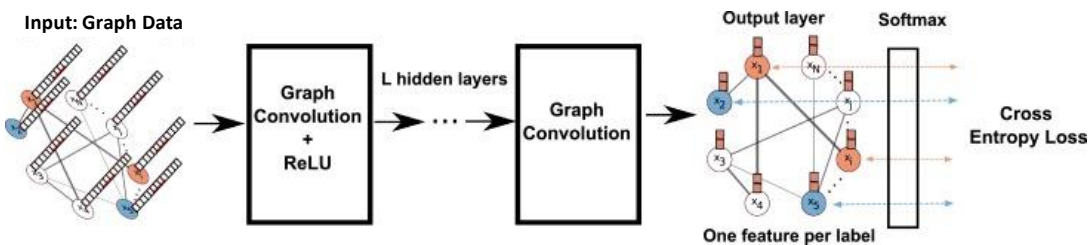
Speech Recognition



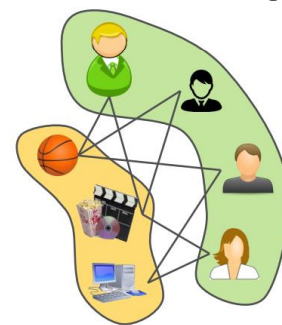
Recommendation



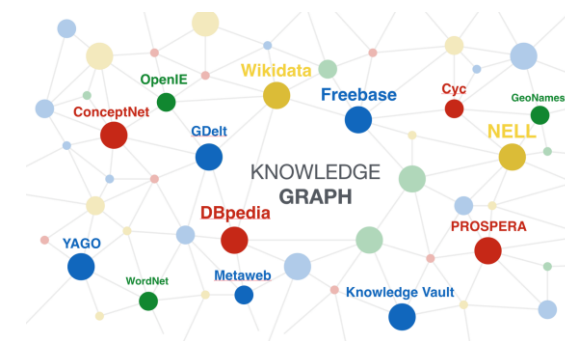
Question Answering



Graph Neural Networks



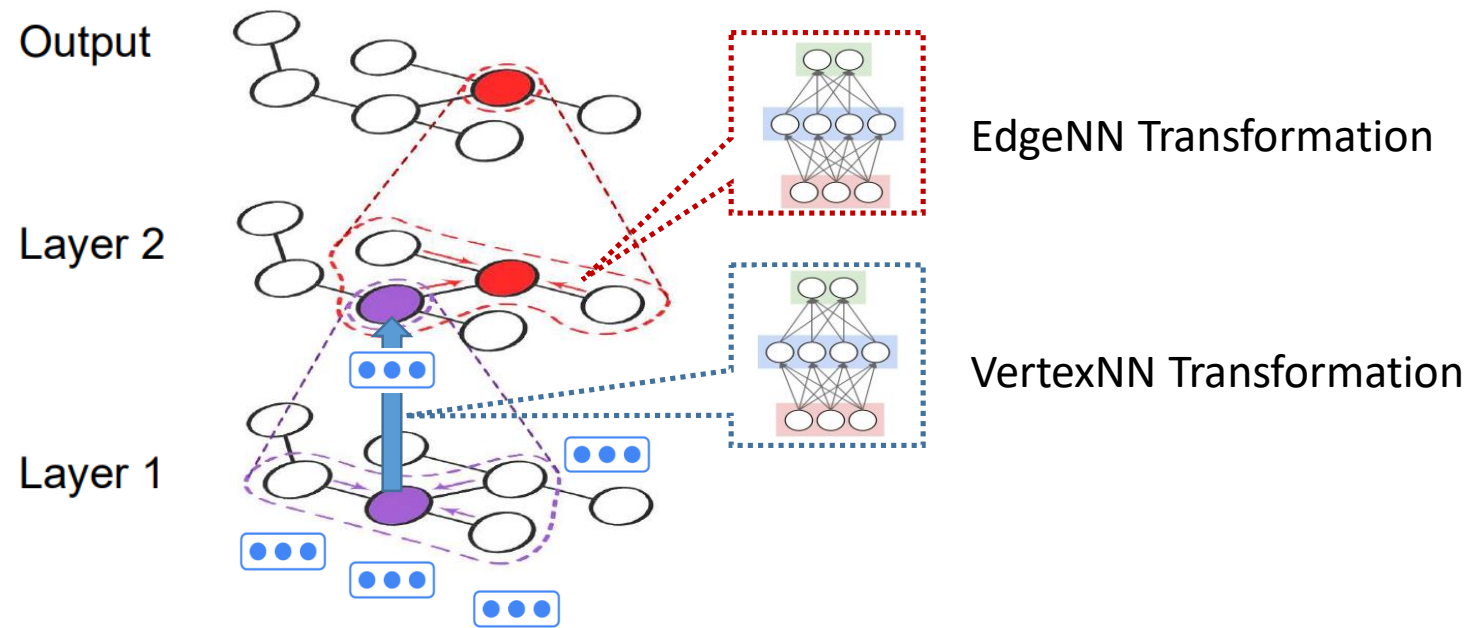
User-Item Graph



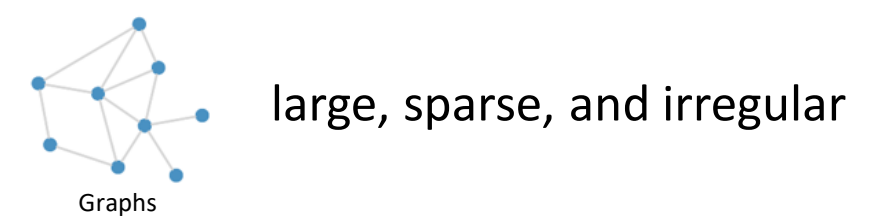
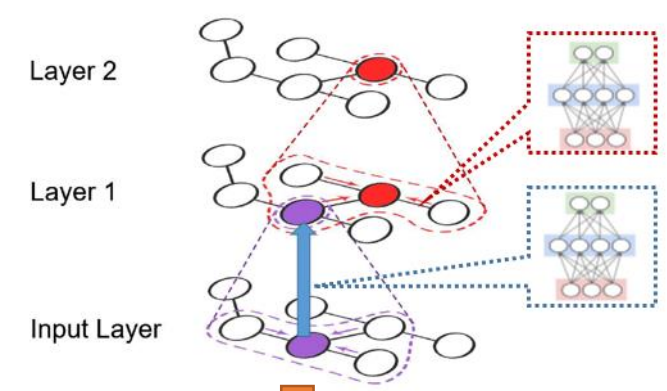
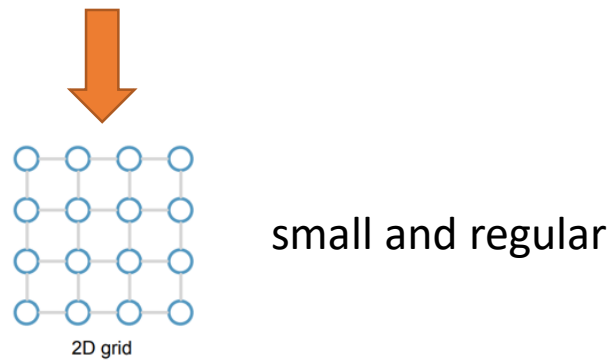
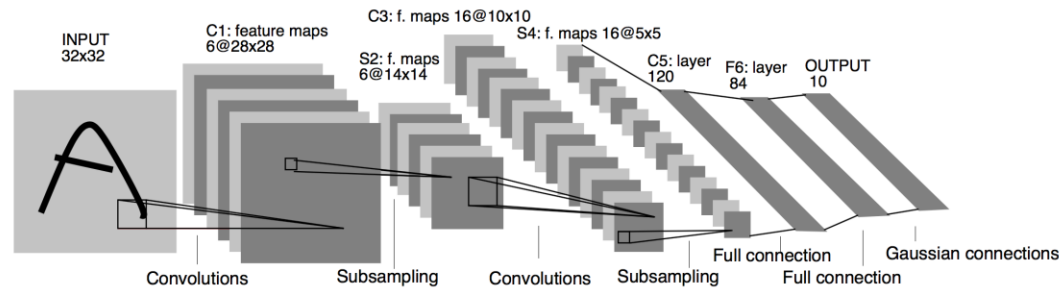
Knowledge Graph

Graph Neural Networks (GNN)

- Information propagation via *Graph*
- Information transformation via *Neural Networks*



Challenges in Processing GNNs on GPU



Existing Systems are Insufficient

Deep Learning Systems

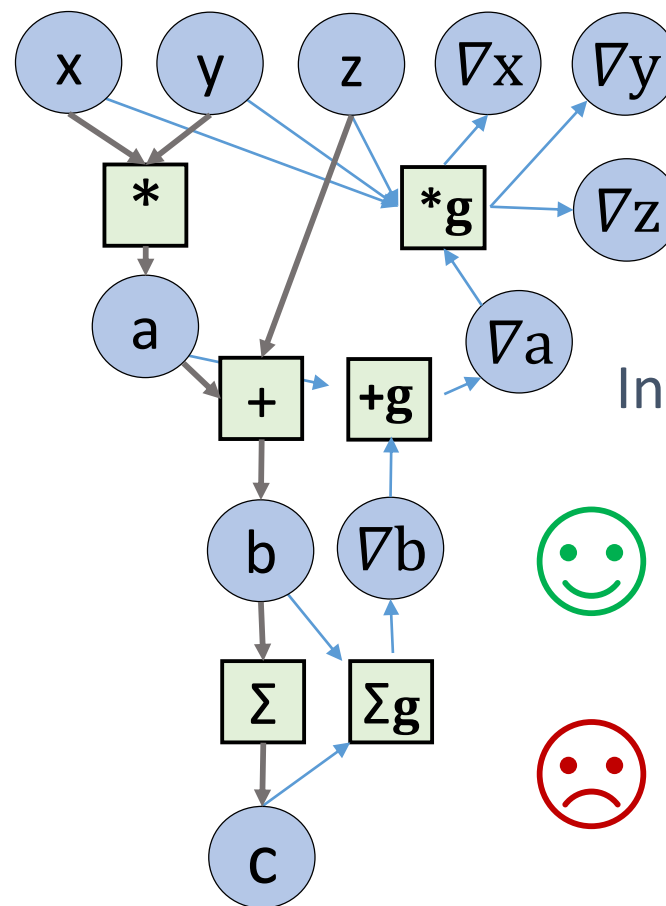


```
x = tf.placeholder(tf.float32)
y = tf.placeholder(tf.float32)
z = tf.placeholder(tf.float32)

a = x * y
b = a + z
c = tf.reduce_sum(b)

grad_x, grad_y, grad_z = tf.gradients(c, [x,y,z])

with tf.Session() as sess:
    sess.run([grad_z], feed_dict=values)
```



Dataflow Graph
as

Intermediate Representation



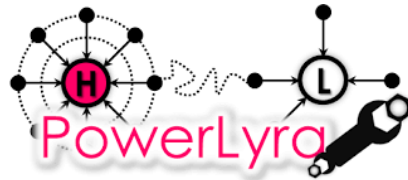
easy to express NNs
efficient for grid structures



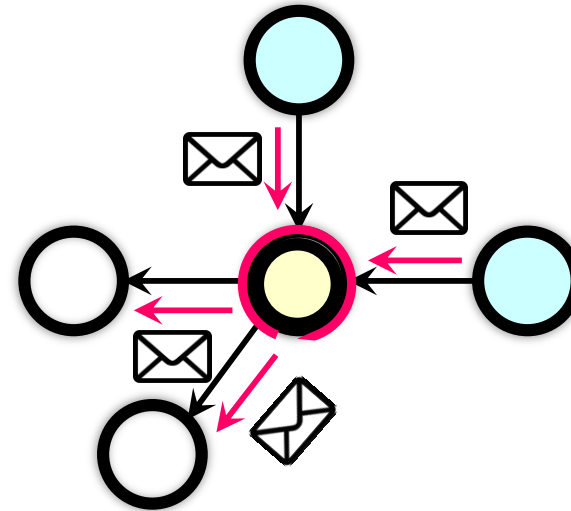
hard to express graph ops
hard to handle large graphs

Existing Systems are Insufficient

Graph Computing Systems



```
Gather( $D_u$ ,  $D_{(u,v)}$ ,  $D_v$ ):  
    return  $D_v.rank / \#outNbrs(v)$   
Sum( $a$ ,  $b$ ):  
    return  $a + b$   
Apply( $D_u$ ,  $acc$ ):  
     $r_{new} = 0.15 + 0.85 * acc$   
     $D_u.delta = (r_{new} - D_u.rank) / \#outNbrs(u)$   
     $D_u.rank = r_{new}$   
Scatter( $D_u$ ,  $D_{(u,v)}$ ,  $D_v$ ):  
    if ( $|D_u.delta| > \epsilon$ ) Activate( $v$ )  
    return  $delta$ 
```



Vertex Programming
e.g.: GAS



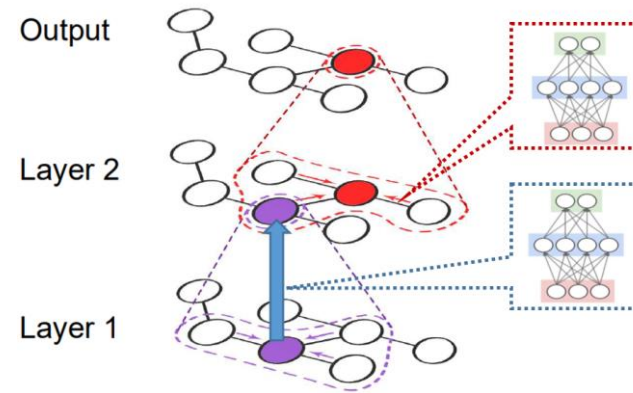
easy to program graph apps
graph-aware optimizations
scale to trillion edges



hard to express NNs (e.g., no backprop.)
insufficient NN execution (e.g., vertex-by-vertex)

We propose: NeuGraph

- Bridge graph and dataflow models to support *efficient* and *scalable* GNN processing



NeuGraph



easy to express NNs
efficient for grid structures

Deep Learning Systems



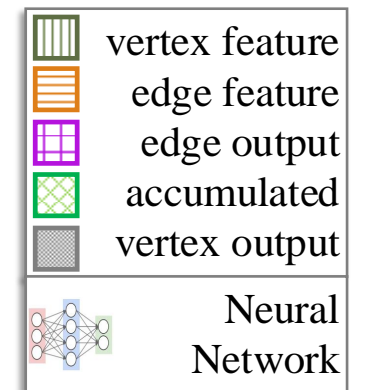
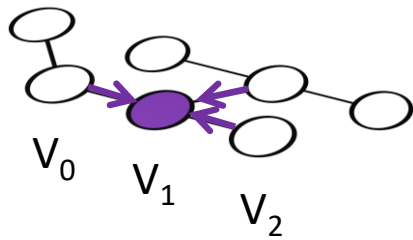
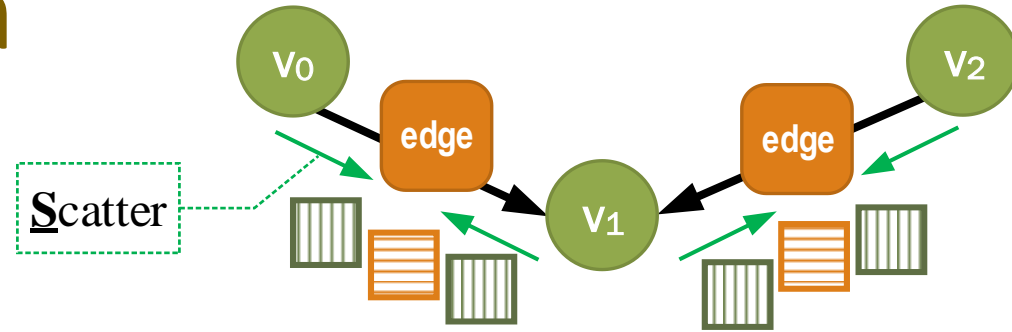
easy to program graph apps
graph-aware optimizations
scale to trillion edges

Graph Systems

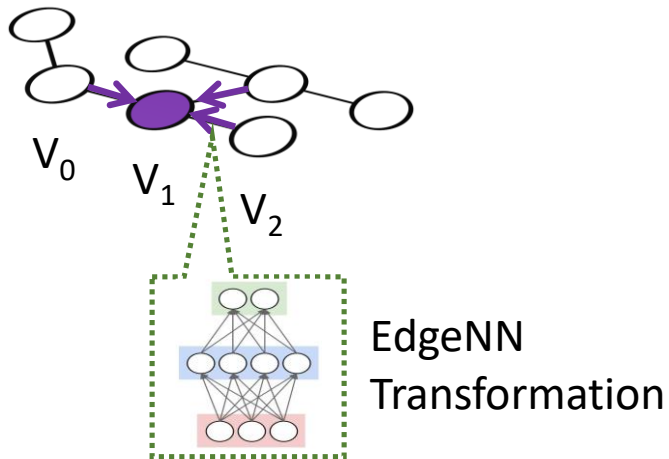
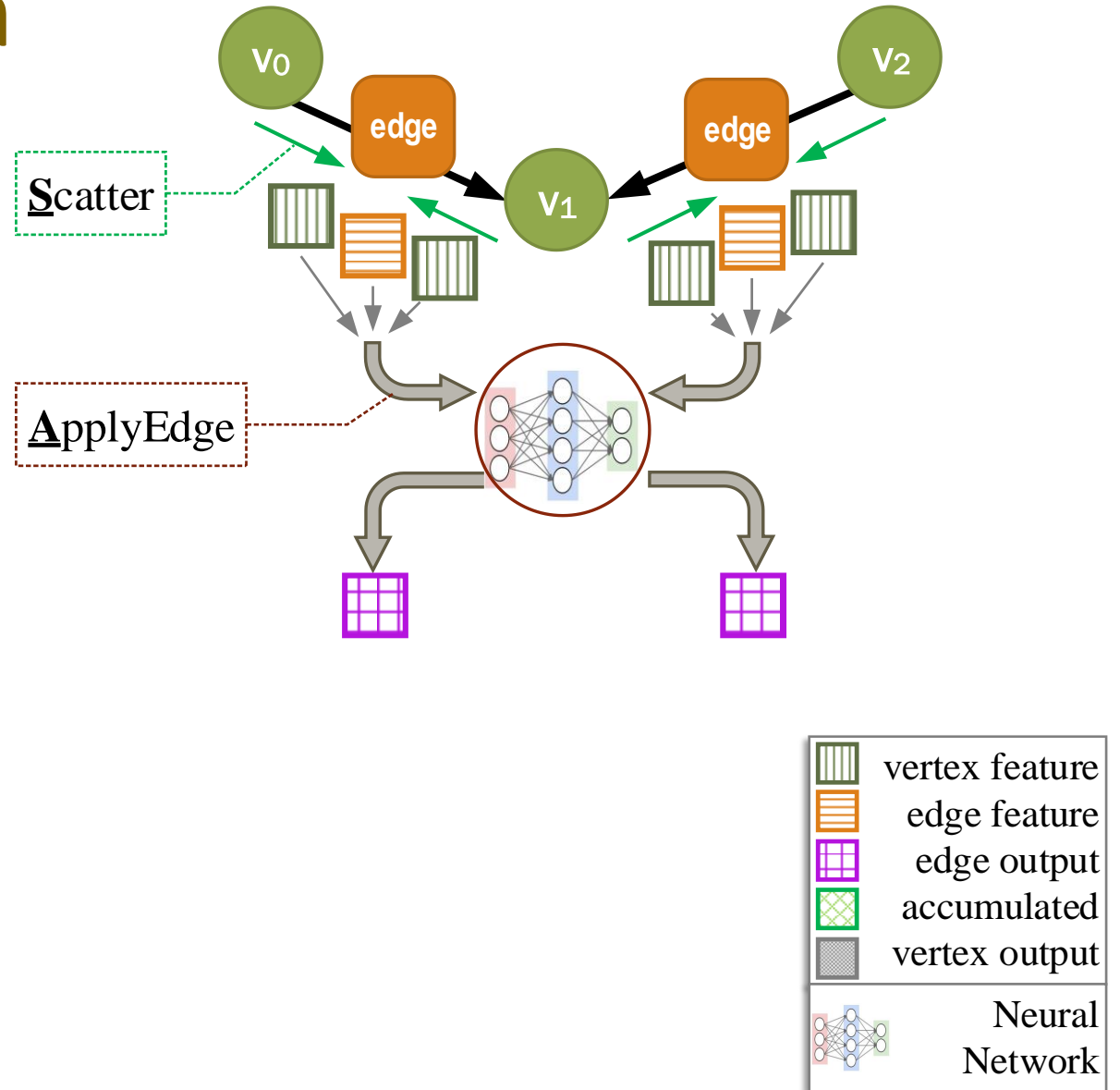
NeuGraph

- Bridge graph and dataflow models to support *efficient* and *scalable* GNN processing
- **Key techniques**
 - SAGA-NN model for graph-based neural networks
 - > *programming GNN apps*
 - Chunk-based dataflow graph translation & streaming processing out of GPU core
 - > *processing graphs larger than GPU memory*
 - Highly-optimized graph propagation operators
 - Chain-based parallel streaming
 - > *efficient multi-GPU parallel execution*
- **Performance**
 - Outperform state-of-the-art frameworks (e.g., TensorFlow and DGL) on small graphs
 - Scale to large real-world graphs with GPUs

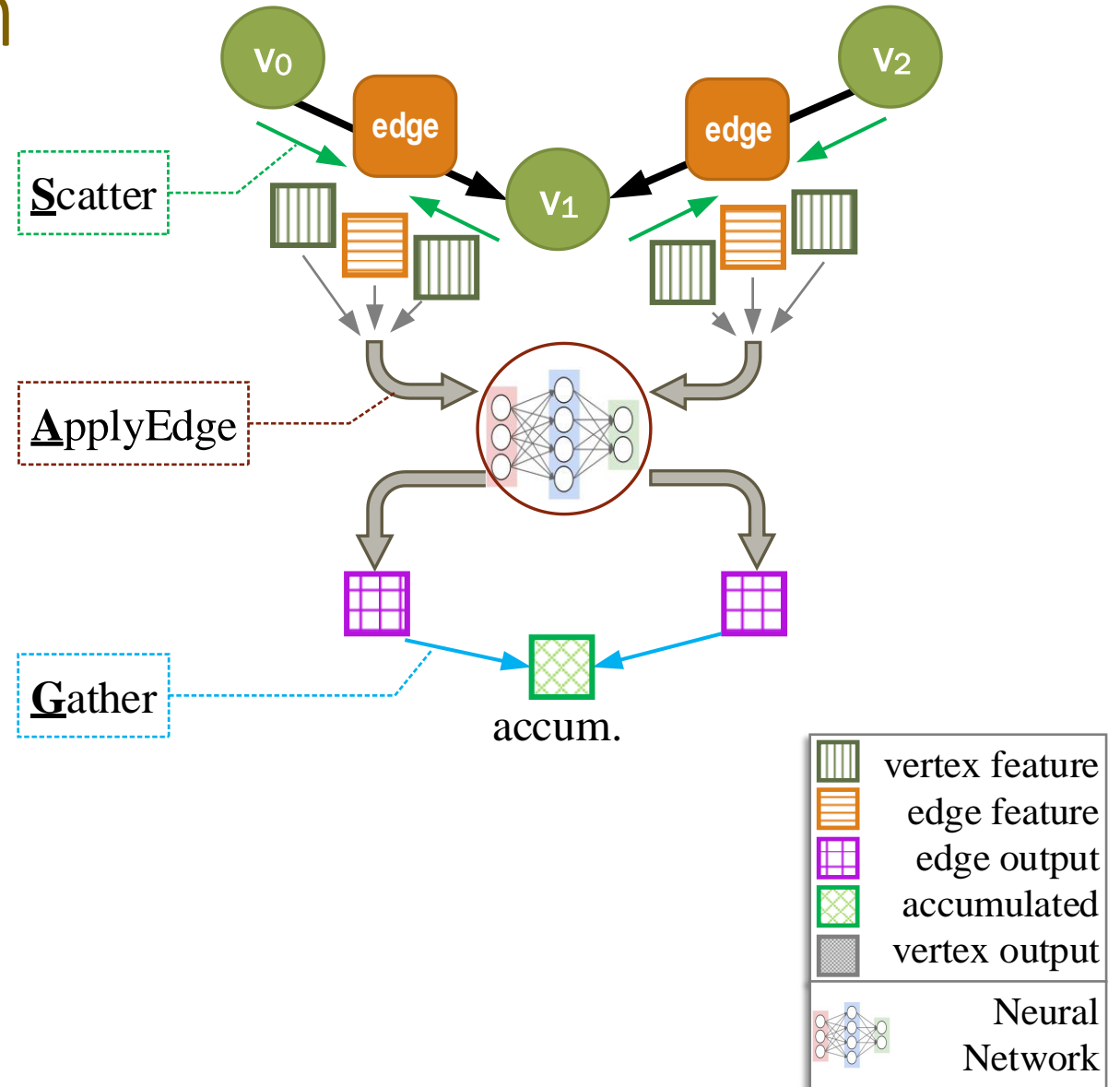
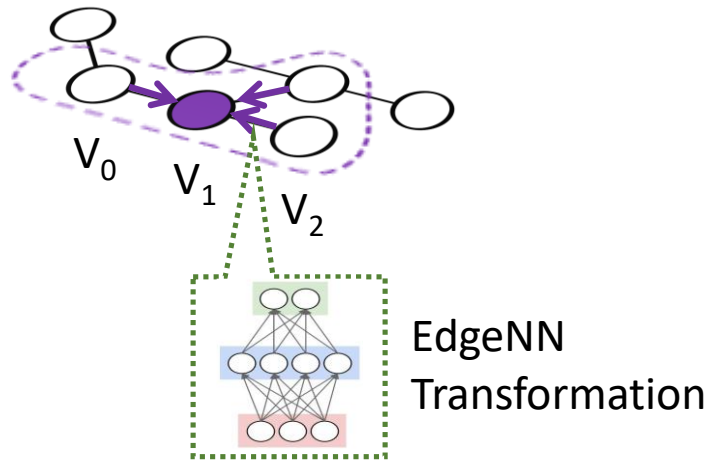
SAGA-NN Abstraction



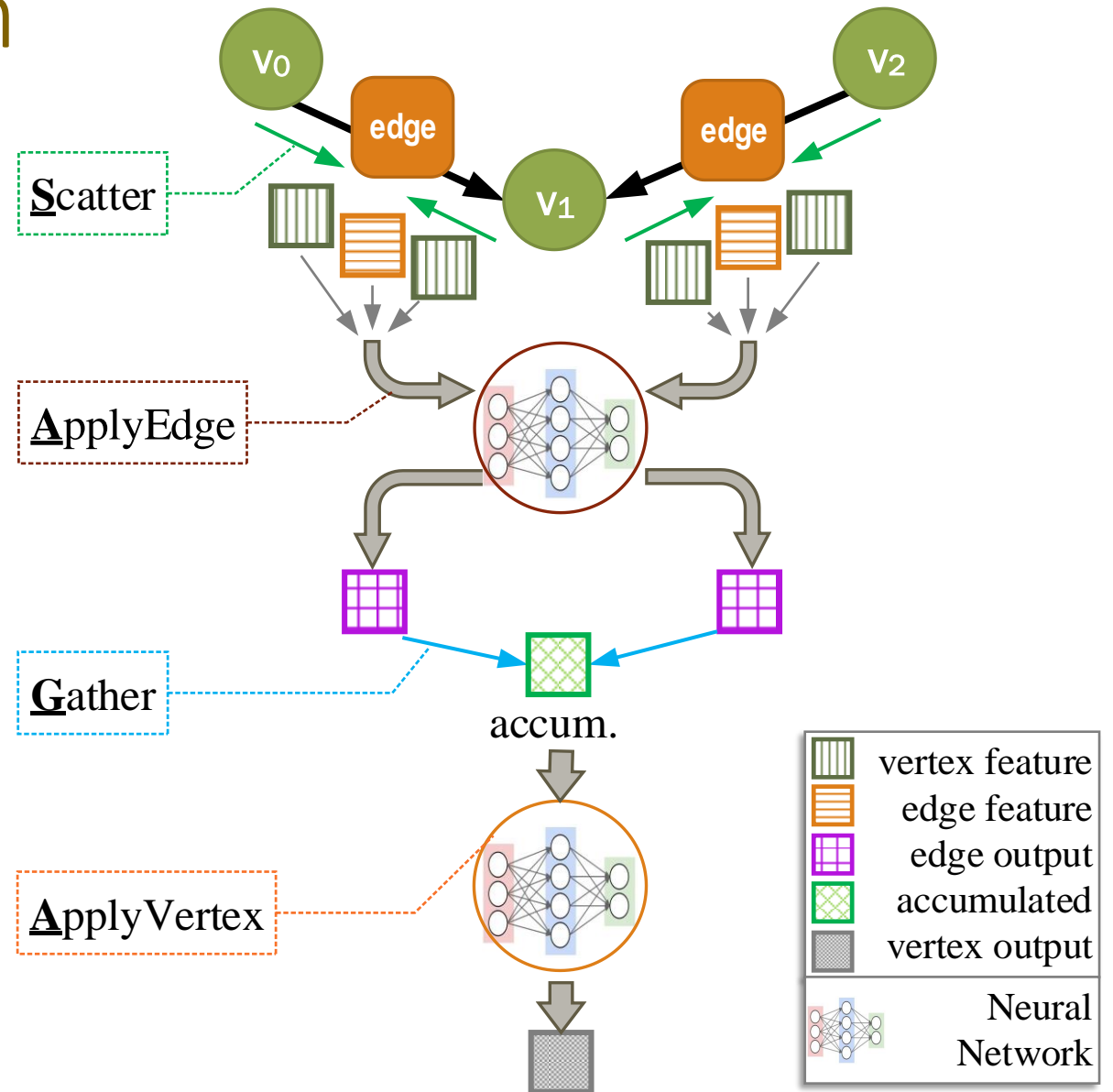
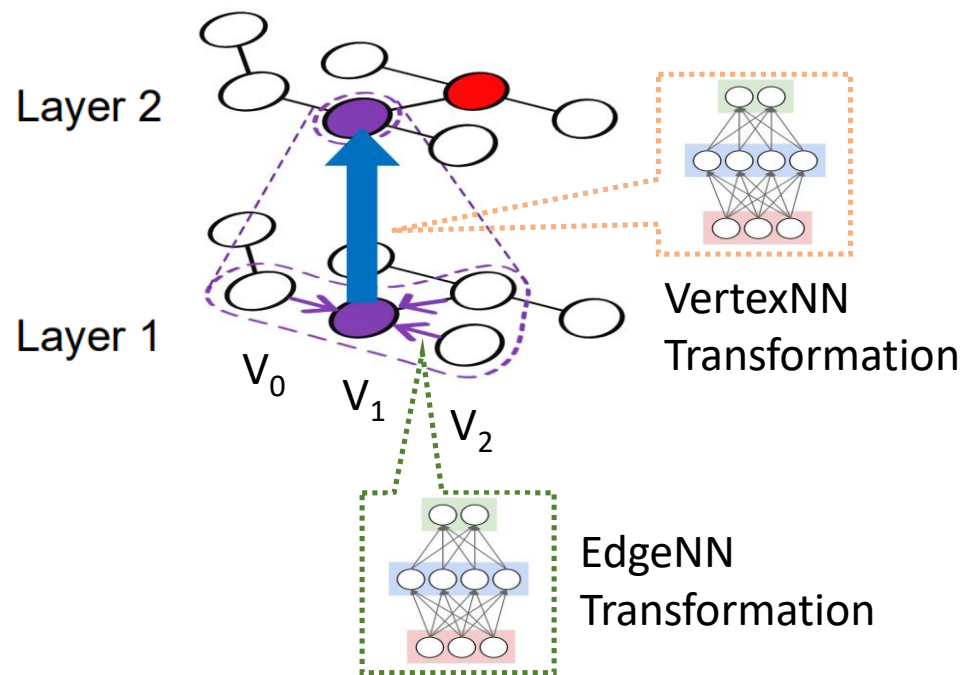
SAGA-NN Abstraction



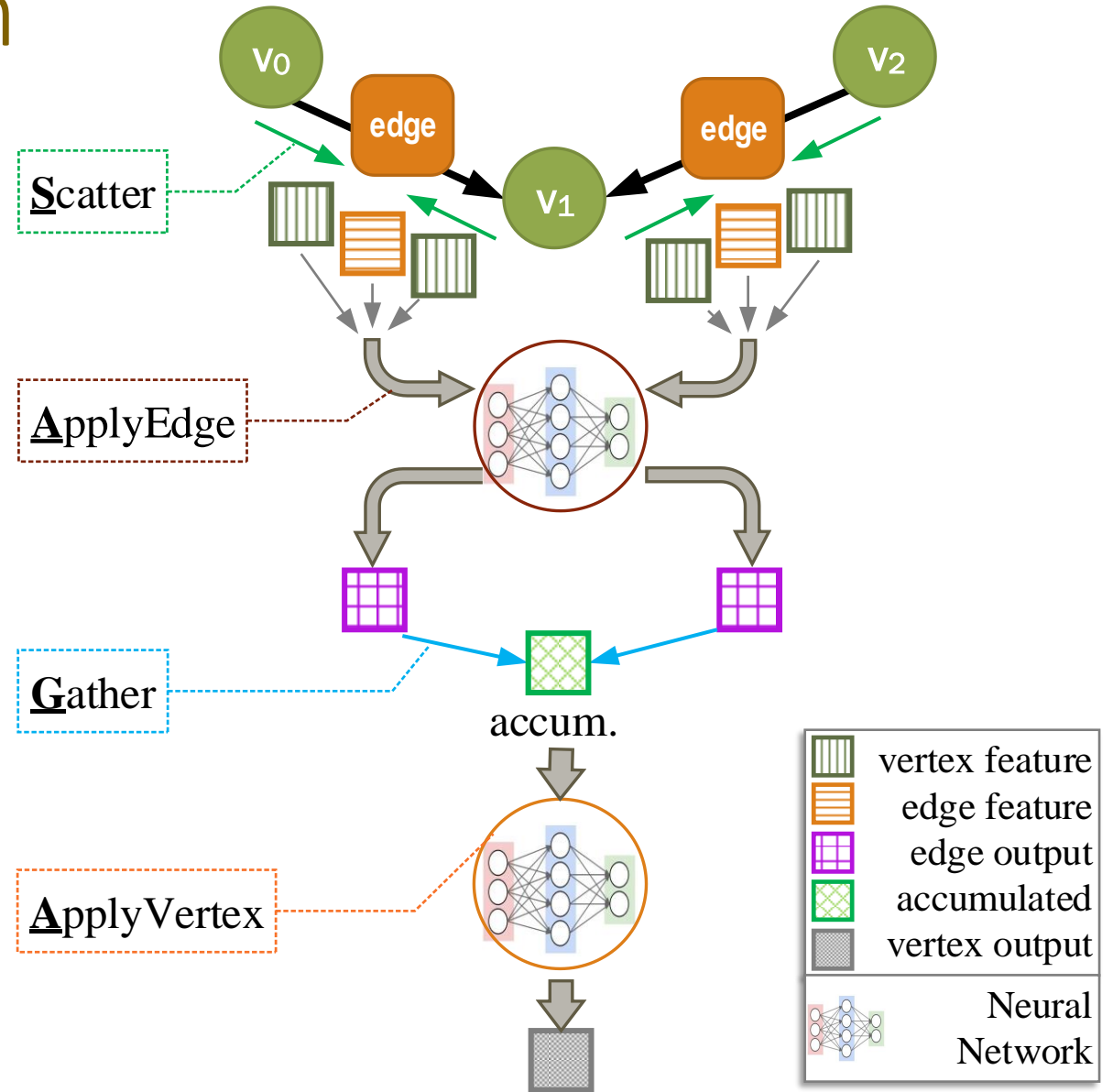
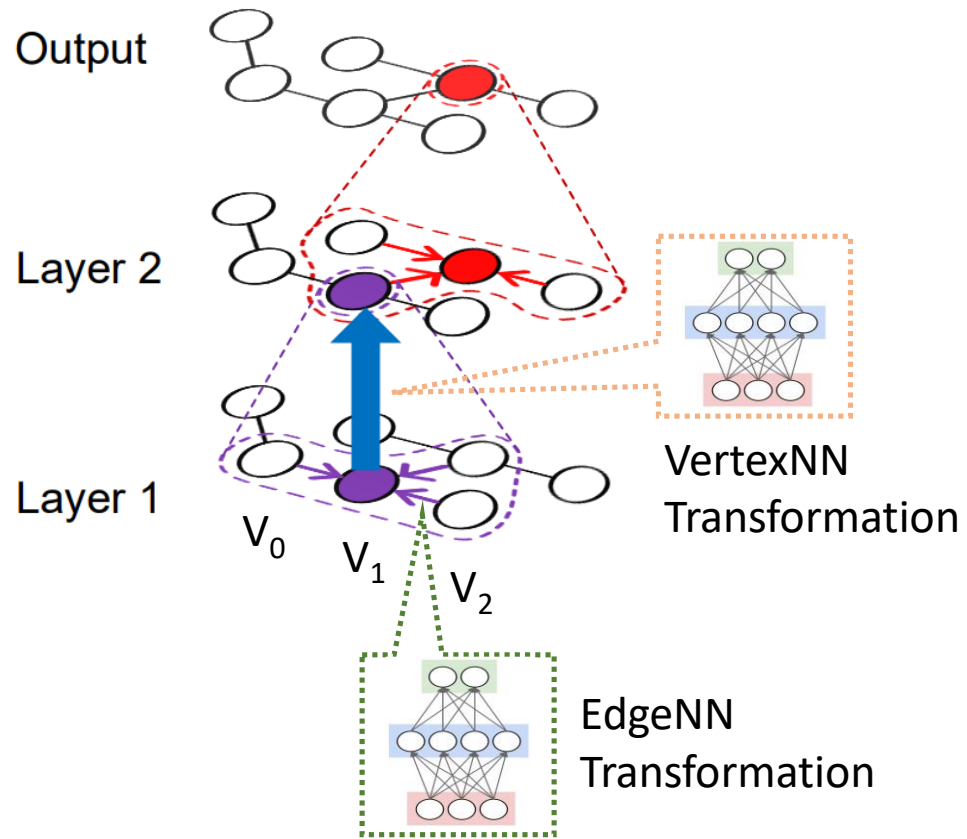
SAGA-NN Abstraction



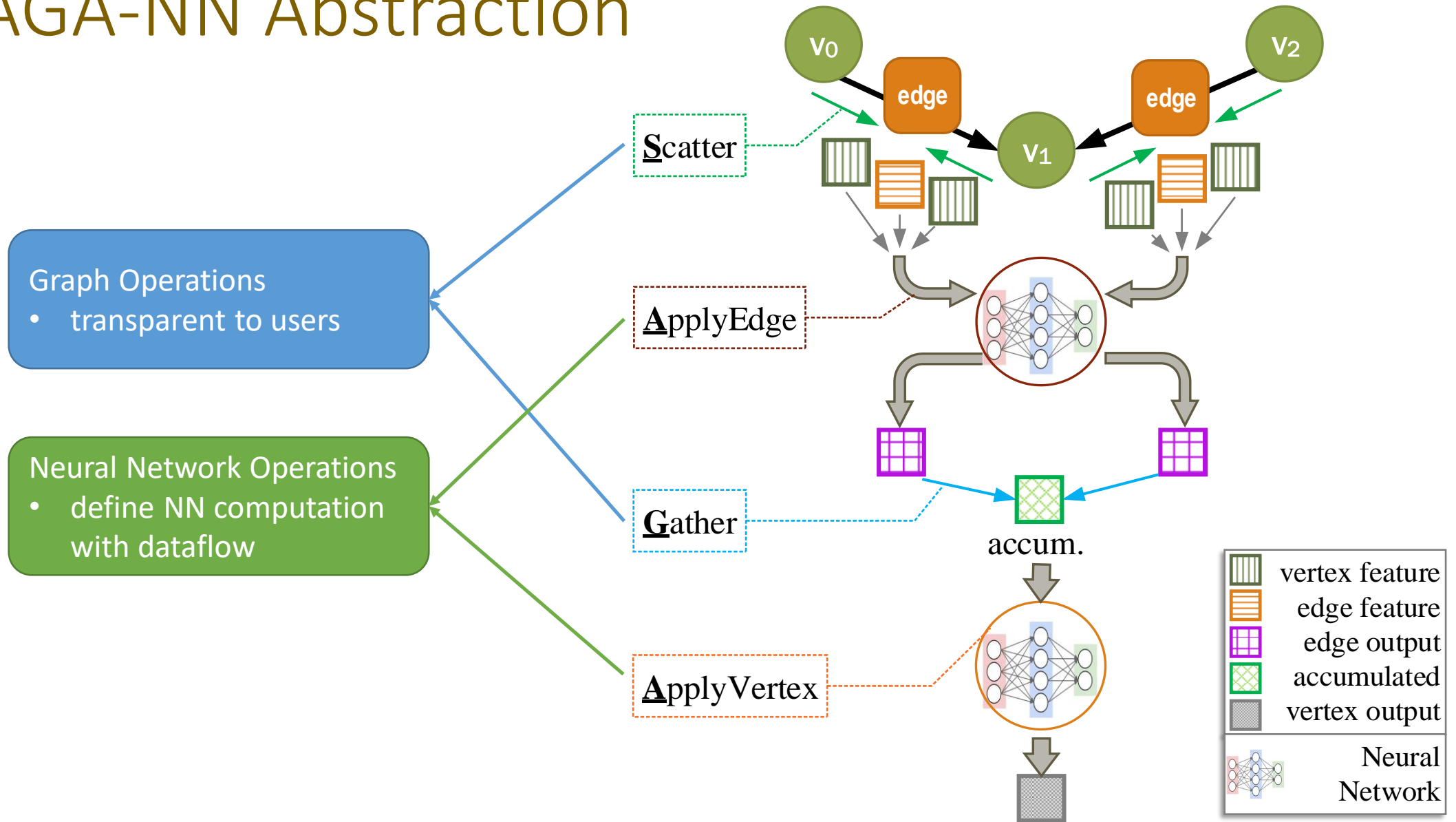
SAGA-NN Abstraction



SAGA-NN Abstraction



SAGA-NN Abstraction



Example: Graph Convolutional Network (GCN)

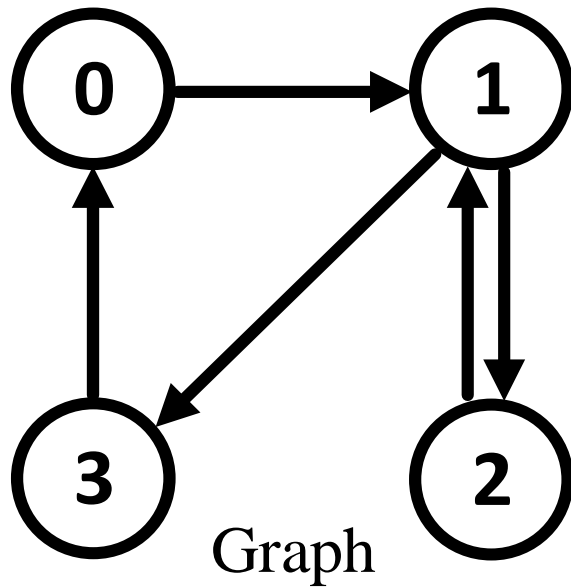
$$\mathbf{h}_u^{l+1} = \text{ReLU} \left(W^l \otimes \left(\sum_{v \rightarrow u} e_{vu} \odot \mathbf{h}_v^l \right) + B^l \right)$$

The diagram illustrates the computation of the next layer's node representation \mathbf{h}_u^{l+1} . The equation is annotated with three callouts:

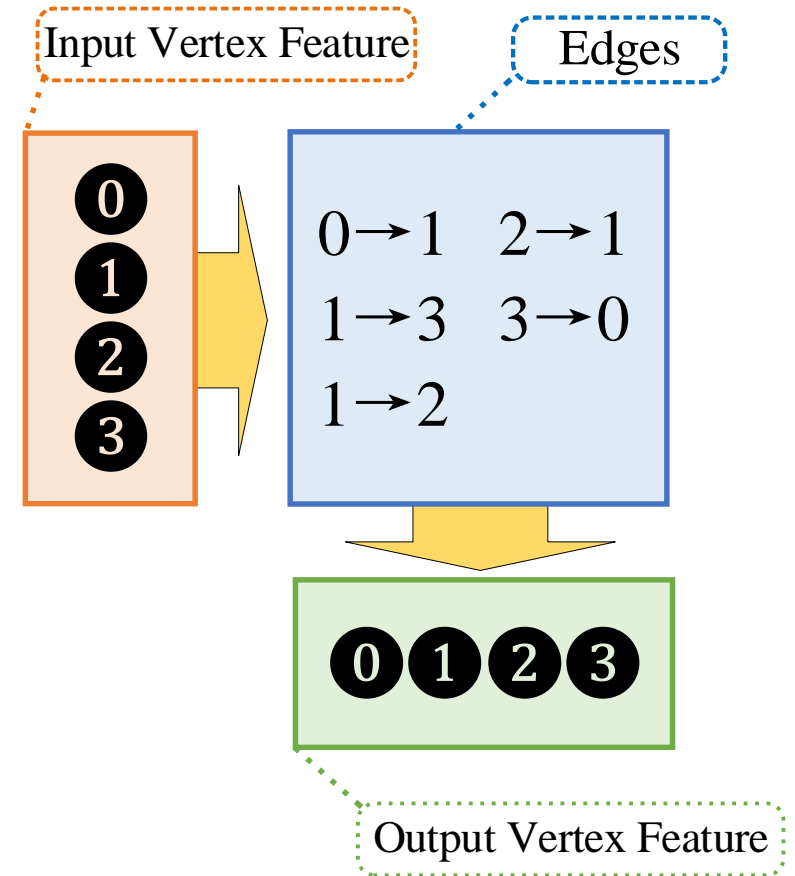
- Gather** (green oval): Points to the summation $\sum_{v \rightarrow u} e_{vu} \odot \mathbf{h}_v^l$, indicating the collection of neighbor node features.
- Scatter ApplyEdge** (blue oval): Points to the element-wise multiplication \odot between the edge weights e_{vu} and the neighbor features \mathbf{h}_v^l .
- ApplyVertex** (orange oval): Points to the entire expression $W^l \otimes \left(\sum_{v \rightarrow u} e_{vu} \odot \mathbf{h}_v^l \right) + B^l$, representing the application of the weight matrix and bias to the gathered and processed neighbor features.

```
1 class GCNLayer(GNNLayer):
2     def __init__(self):
3         self.gather_accumulator = "sum"
4         ... # init variables
5     def _apply_vertex(self, vertex, accum):
6         ret = tf.matmul(accum, self.vars['weight'])
7         ret = ret + self.vars['bias']
8         return tf.nn.relu(ret)
9     def _apply_edge(self, edge):
10        return edge.src * edge.data
11 class GCNApp(GNNModel):
12     def __init__(self):
13         ... # init configs, optimizer, loss, et.al.
14     def _build(self):
15         self.layers.append(GCNLayer(...)) # params
16         self.layers.append(GCNLayer(...))
17         self.layers.append(GCNLayer(...))
18         self.layers.append(GCNLayer(...))
```

Scaling beyond GPU memory limit

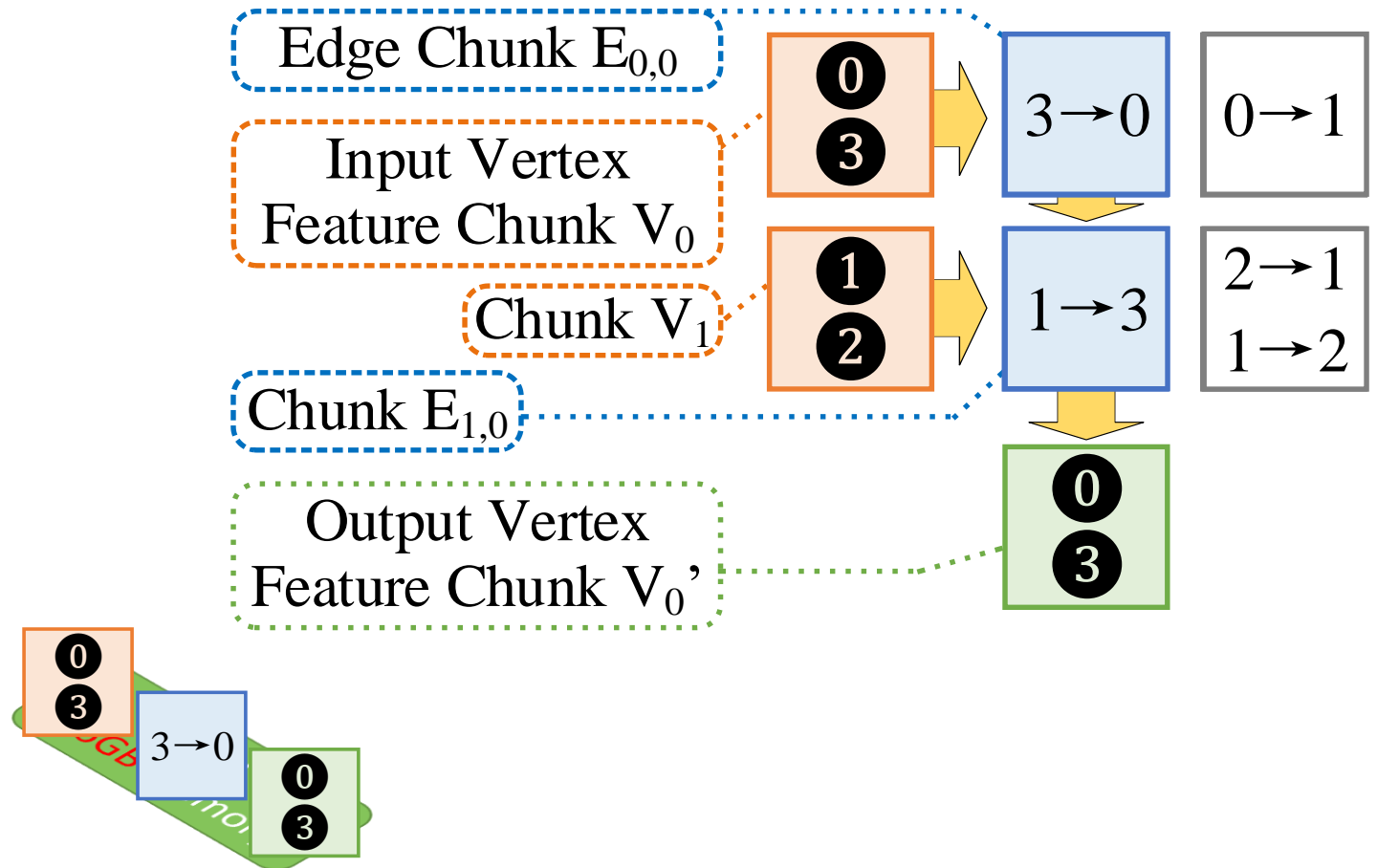
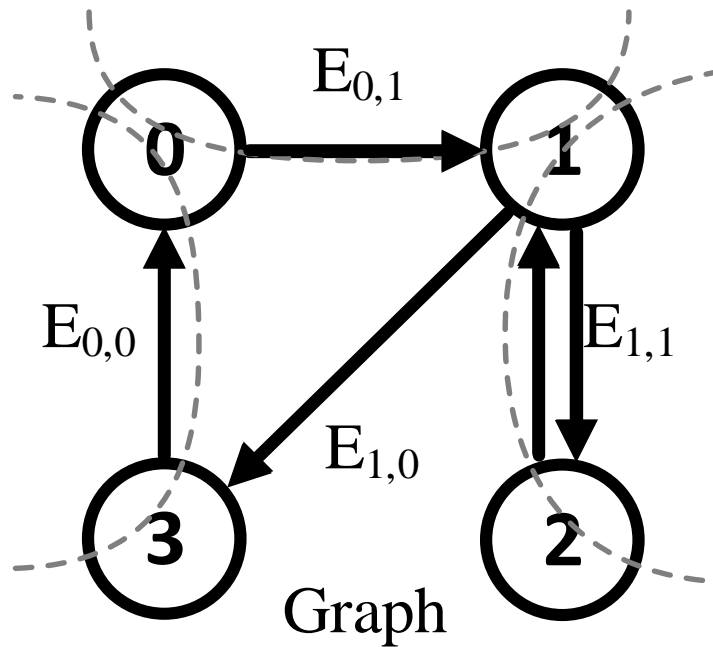


GPU Memory
(16GB)



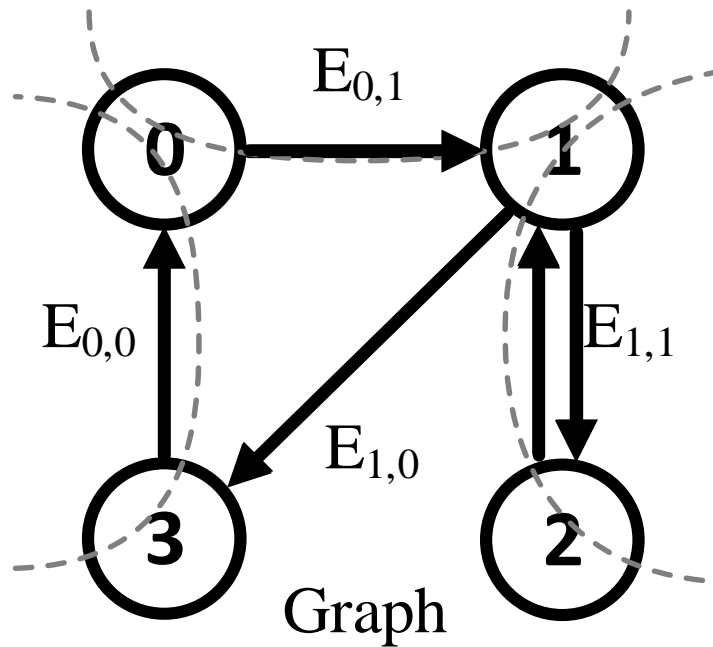
Scaling beyond GPU memory limit

- 2D graph partitioning

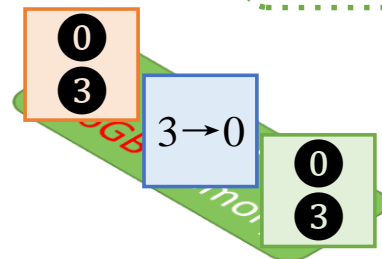
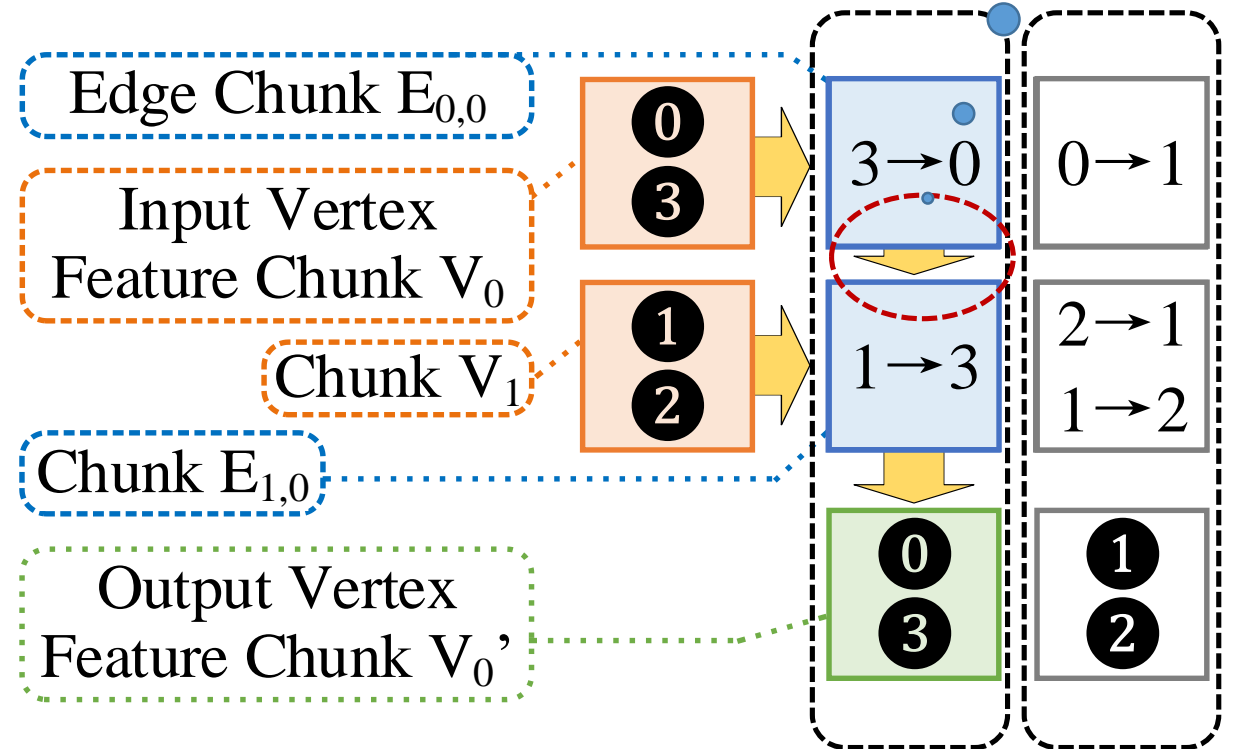


Scaling beyond GPU memory limit

- 2D graph partitioning



Accumulate in GPU memory



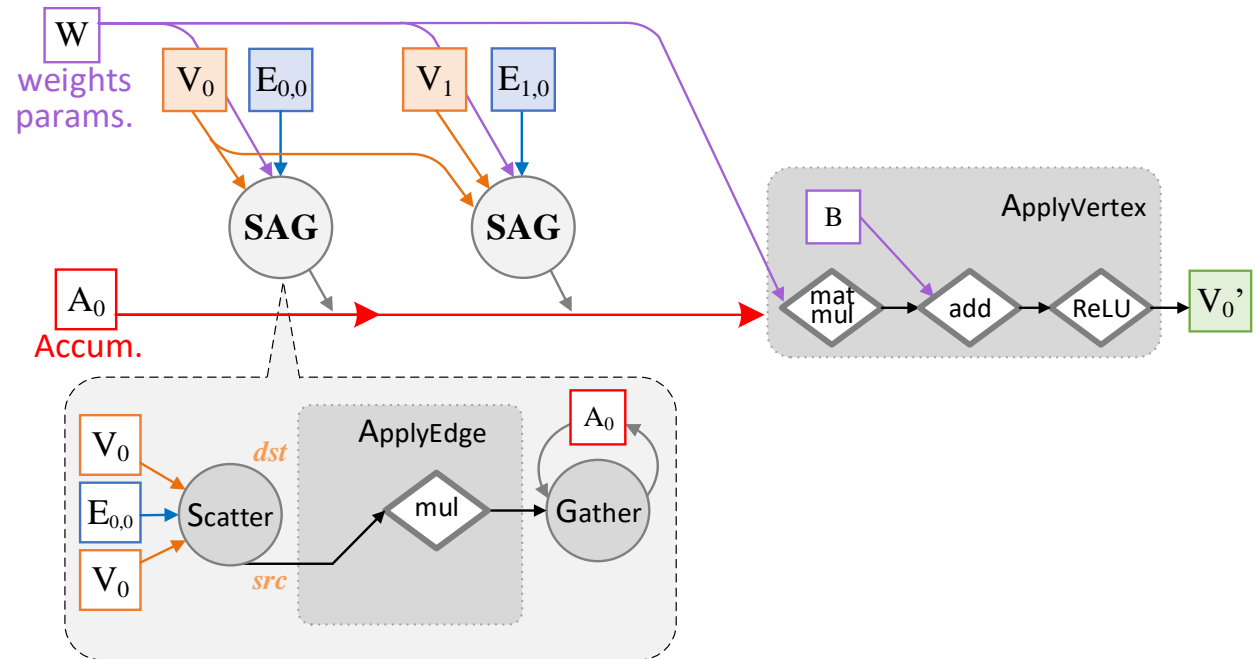
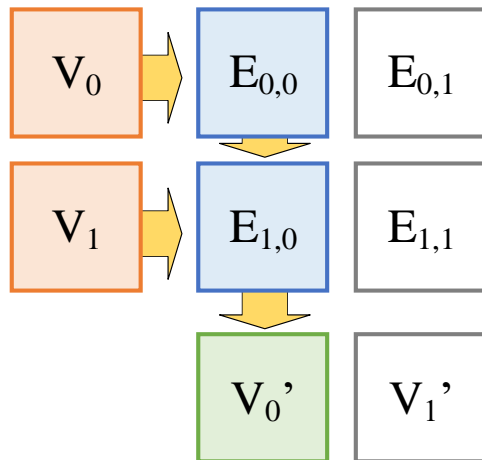
column-by-column
chunk scheduling

Chunk-based Dataflow Translation: GCN

```

1 class GCNLayer(GNNLayer):
2     def _init(self):
3         self.gather_accumulator = "sum"
4         ... # init variables
5     def _apply_vertex(self, vertex, accum):
6         ret = tf.matmul(accum, self.vars['weight'])
7         ret = ret+self.vars['bias']
8         return tf.nn.relu(ret)
9     def _apply_edge(self, edge):
10        return edge.src * edge.data

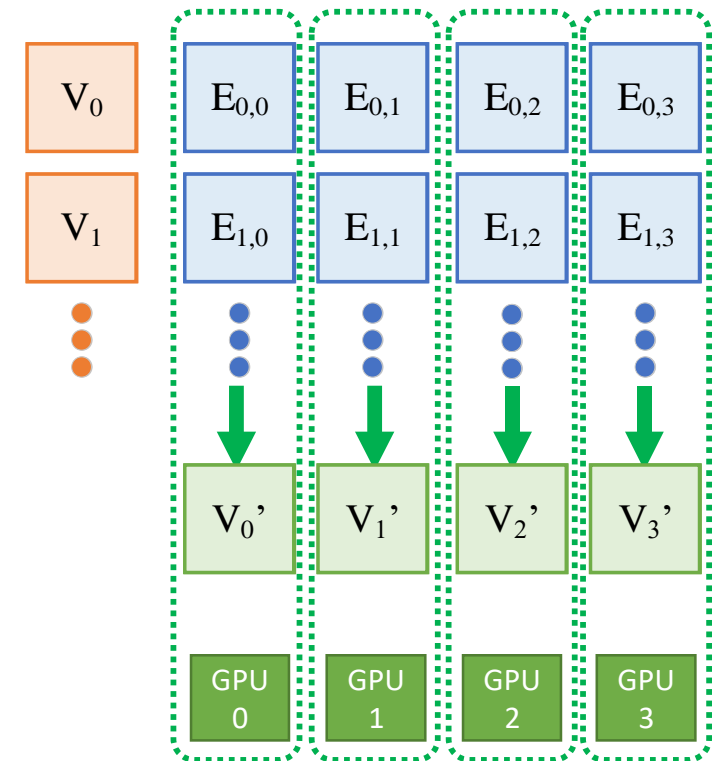
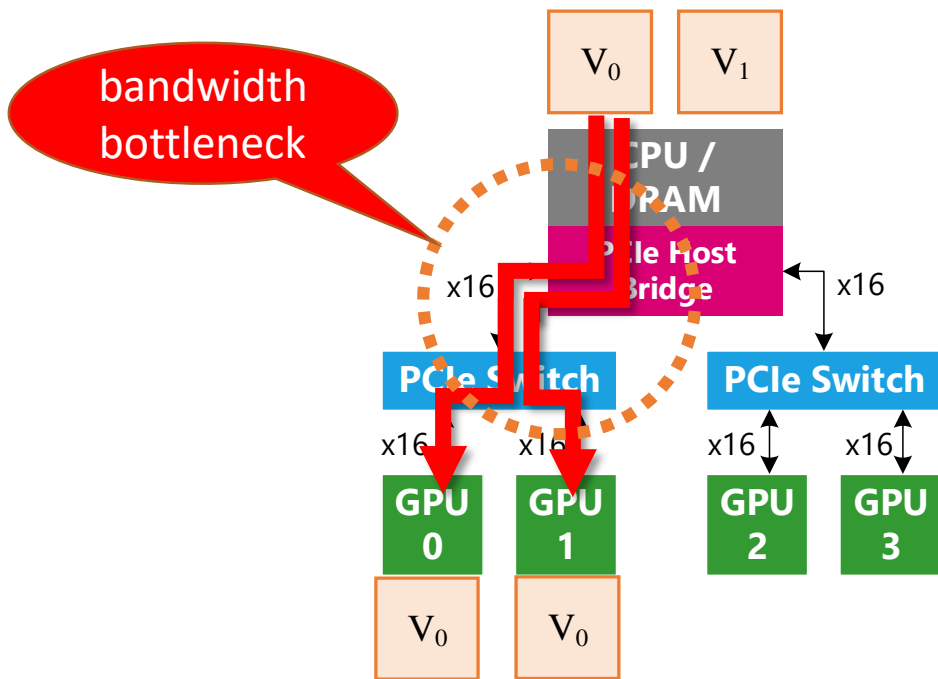
```



Scaling to multi-GPU

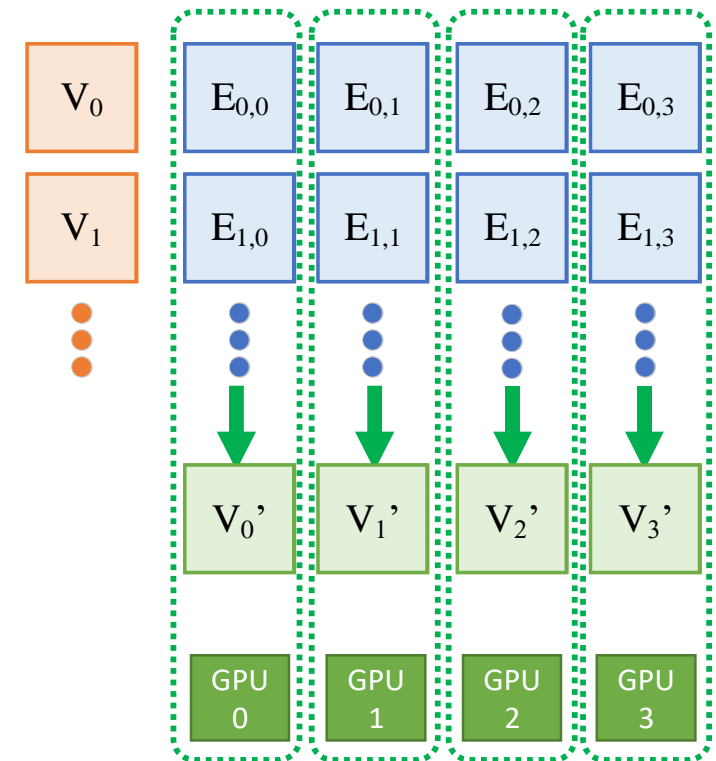
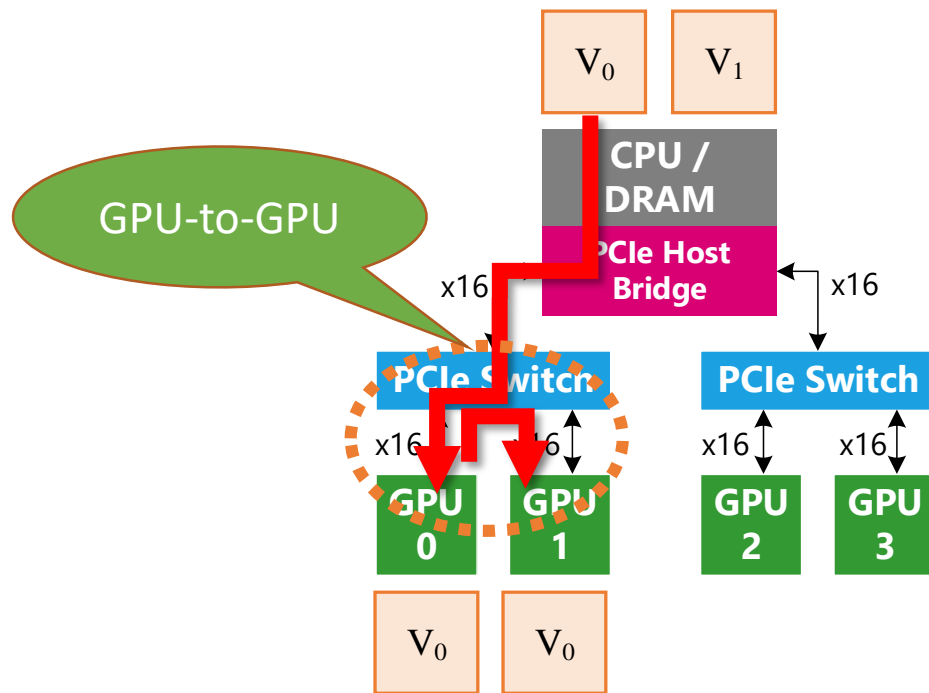
- Naïve solution:

- Approach: each GPU takes a “column” of edge chunks
- Problem: redundant data transfer through shared upper links



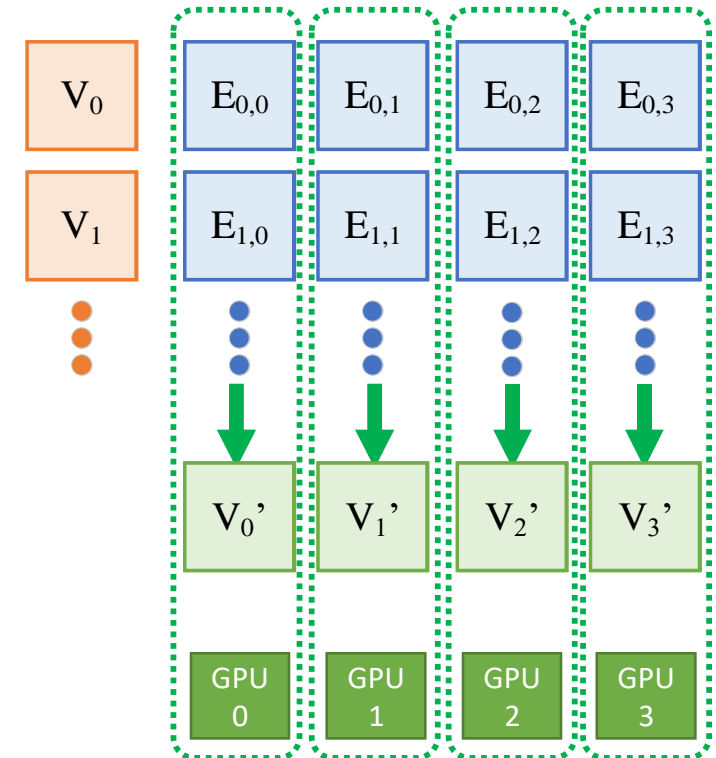
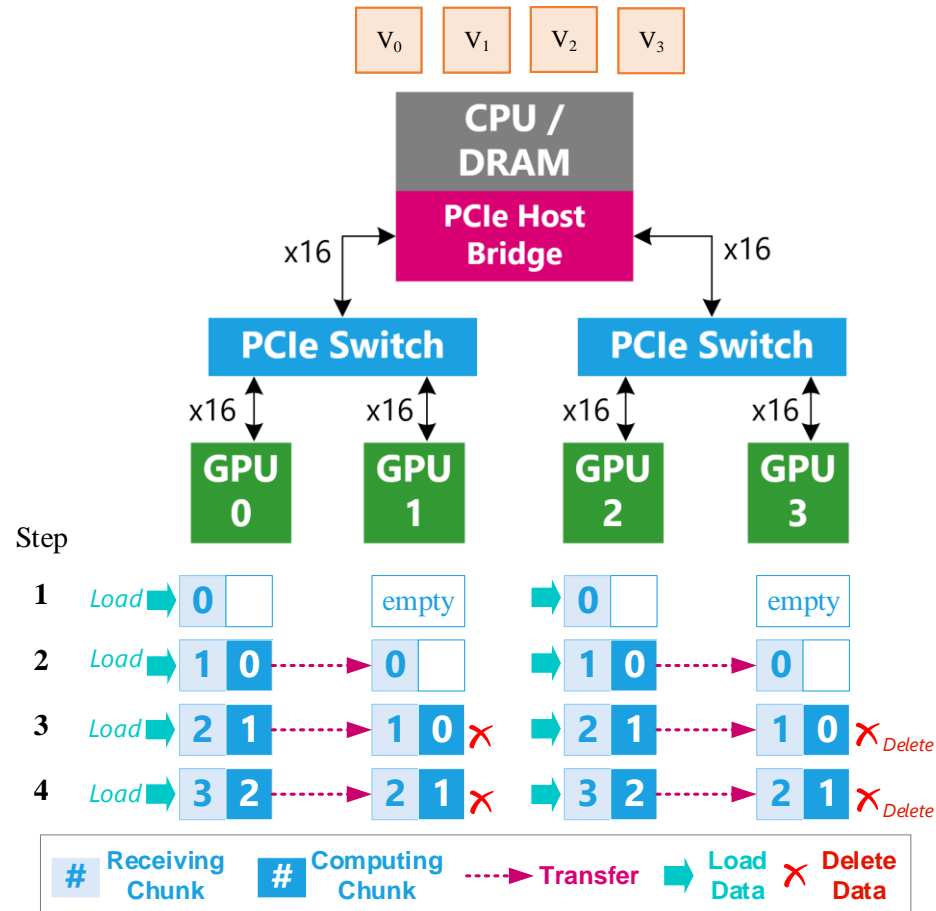
Scaling to multi-GPU

- Observation: link to neighbor GPU is empty
- Approach: directly load from neighbor GPU



Chain-based Parallel Streaming

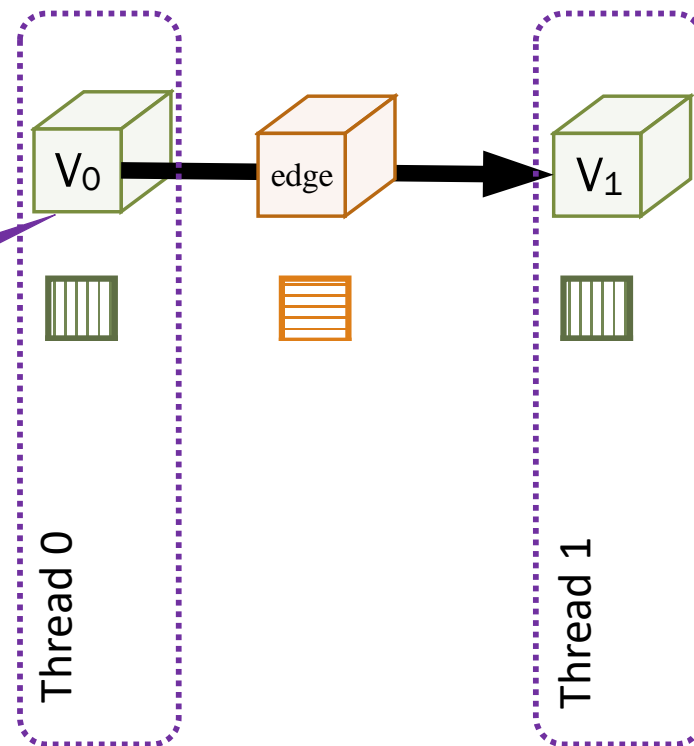
- Each GPU can consume all the chunks in chain for a PCIe switch



High performance graph propagation kernels

- Traditional multithread parallelization

- PageRank, SSSP, CC, etc: one-dimension scalar value on vertices
- A thread takes a vertex/edge



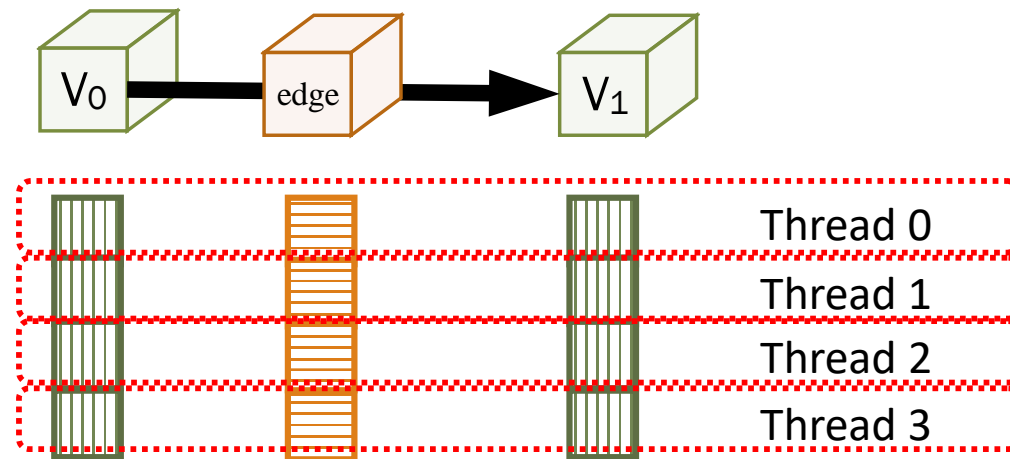
diff ops. on diff vertices
inefficient in GPU(SIMD)

High performance graph propagation kernels

- Observation

- High-dimension feature
- More-regular operations along feature-dimension

- Solution: exploit feature-dimension parallelization in GNNs



same ops.
efficient in GPU(SIMD)

Experiment Setup

- Typical GNN Applications

- Communication neural network (**CommNet**)
// no computation in ApplyEdge
- Graph convolutional networks (**GCN**)
// element-wise Mul in ApplyEdge
- Gated graph neural networks (**GG-NN**)
// GRU in ApplyVertex

- Testbed

- 2x E5-2690-v4 (14-core with HT)
- 512GB Quad-Channel RAM
- 8x NVIDIA Tesla P100 GPU
- Ubuntu 16.04
- CUDA 9.0 + cuDNN 7

dataset	vertex	edge	feature	label	degree	type
blog	10.3K	668.0K	128	39	65	social
pubmed	19.7K	108.4K	500	3	5	cite
reddit_small	58.2K	1.4M	300	41	25	social
reddit_full	2.4M	705.9M	300	50	292	social
enwiki	3.6M	276.1M	300	12	77	knowledge
amazon	8.6M	231.6M	96	22	27	review

- Comparison

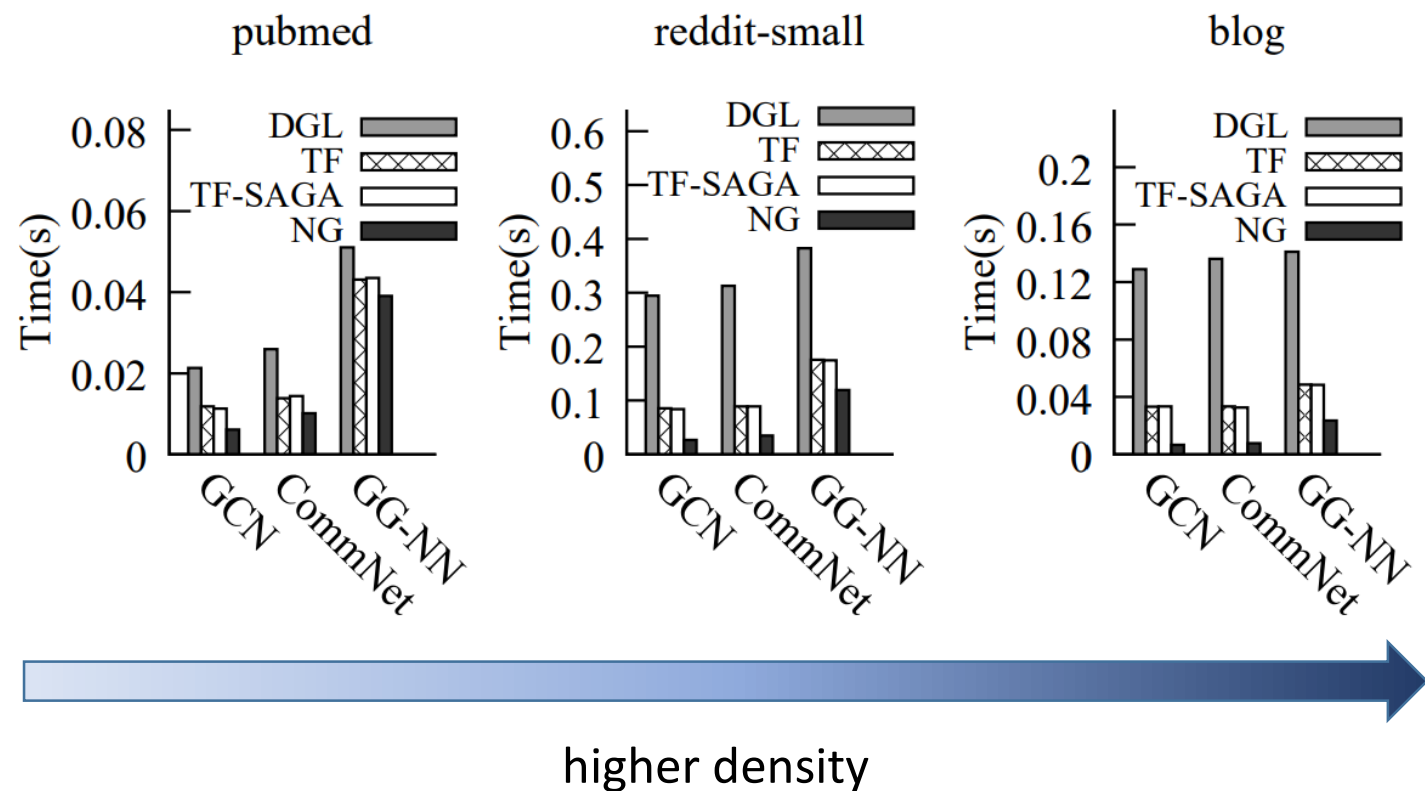
- TF: TensorFlow v1.7 (CPU/GPU)
- DGL: Deep Graph Library v0.1.3 with PyTorch v1.0 backend
- TF-SAGA: SAGA-NN with TensorFlow v1.7 backend
- NG: NeuGraph on top of TensorFlow v1.7

Performance on a Single GPU

Compare with TF, DGL on small graphs

up to 5x speedup vs TF

up to 19x speedup vs DGL



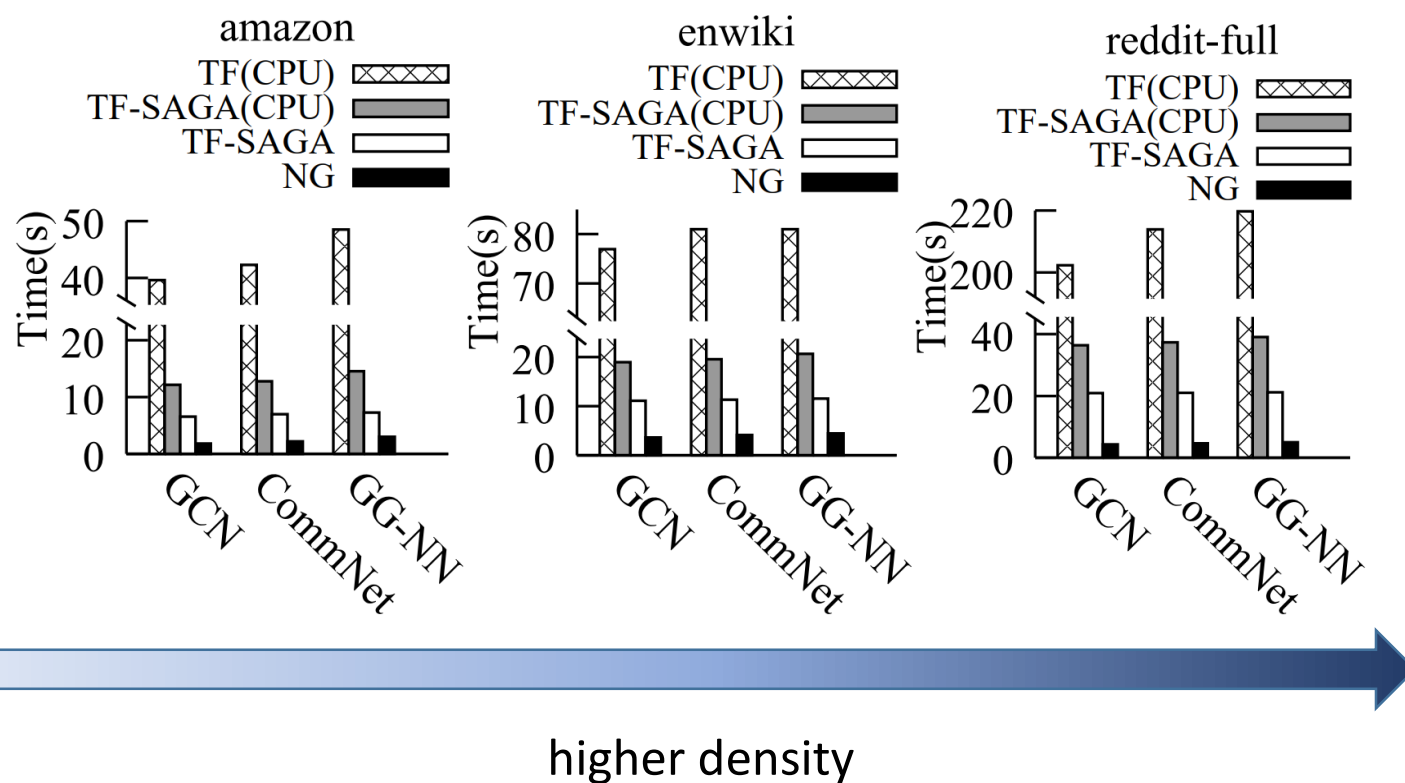
Avg. epoch time; #layer=2, hidden_size=512

Scaling-up on a Single GPU

- Compare with TF(CPU), TF-SAGA(CPU), TF-SAGA on large graphs
- TF and DGL run out-of-memory (OOM)

up to 47x speedup vs TF(CPU)

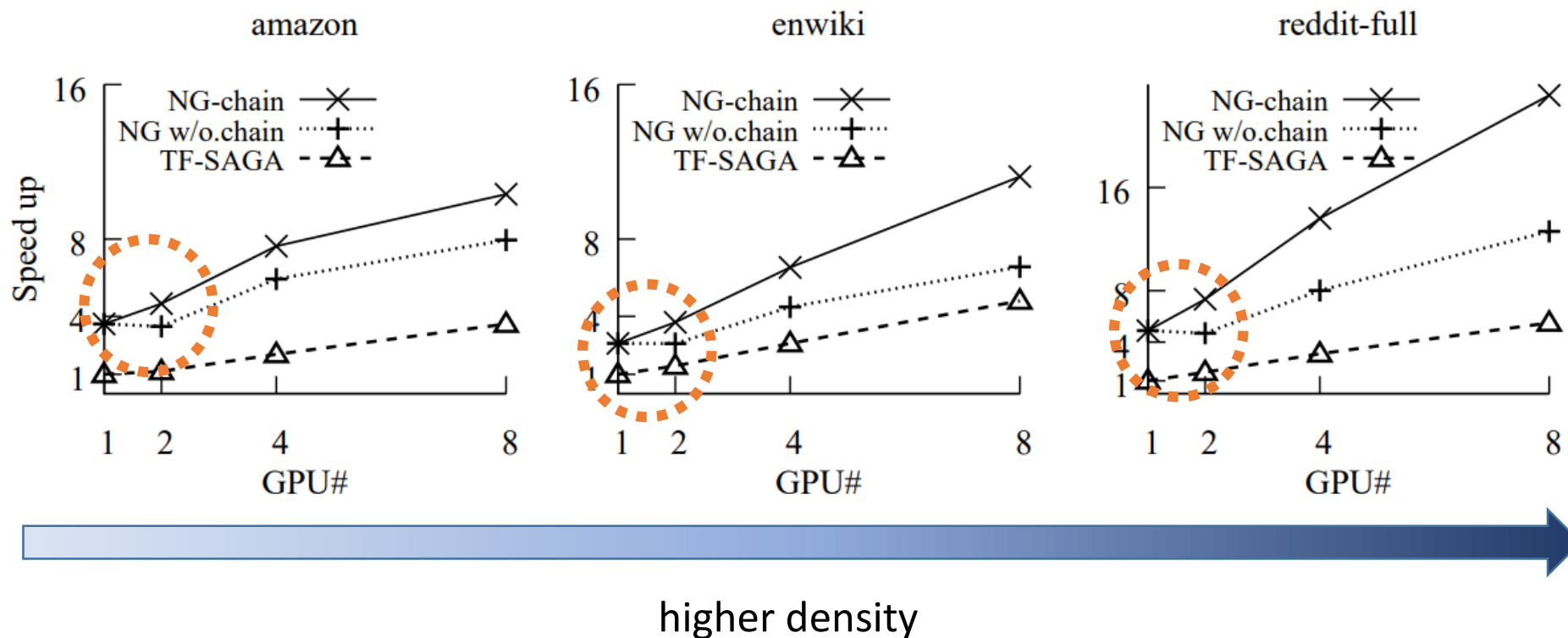
up to 5x speedup vs TF-SAGA



Avg. epoch time; #layer=2, hidden_size=16 (amazon), 32 (enwiki), 64 (reddit-full)

Scaling-out on Multiple GPUs

- GCN on three large graphs on different number of GPUs



Avg. speedup over TF-SAGA-1GPU; #layer=2, hidden_size=16 (amazon), 32 (enwiki), 64 (reddit-full)

Conclusion

NeuGraph: advocate unifying deep learning and graph computing for graph neural networks

- Define a new, flexible SAGA-NN model to express GNNs
- Fuse graph-related optimizations into deep learning frameworks
- Demonstrate NeuGraph achieves efficient and scalable GNN processing

NeuGraph

Thank you!
Q&A

For more information, please visit: <https://aka.ms/neugraph>



Microsoft®
Research