

Scaling Up Graph Neural Network

PRESENTER: Ningyi Liao

SUPERVISOR: Asst. Prof. Siqiang Luo

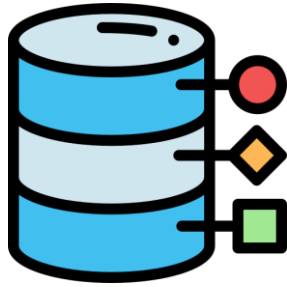


Outline

Introduction



Design Scalable Algorithms



Deploy Scalable Data Structures



Evaluate Scalability Performance

Conclusion and Future Directions

Outline

Introduction

Graph Data and Representation

Graph Neural Network

Challenge and Objectives



Algorithms



Data Structures



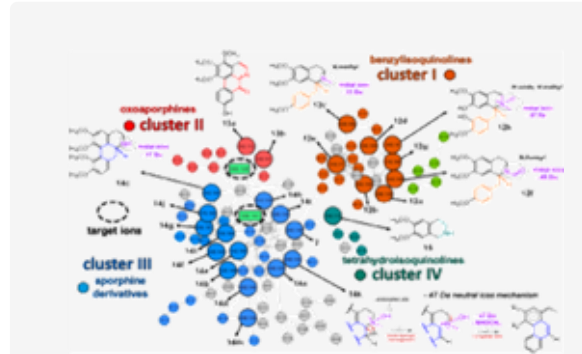
Evaluations

Conclusion and Future Directions

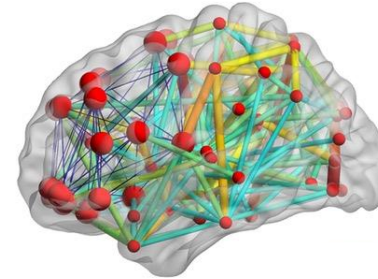
Graph: A Ubiquitous Data Structure



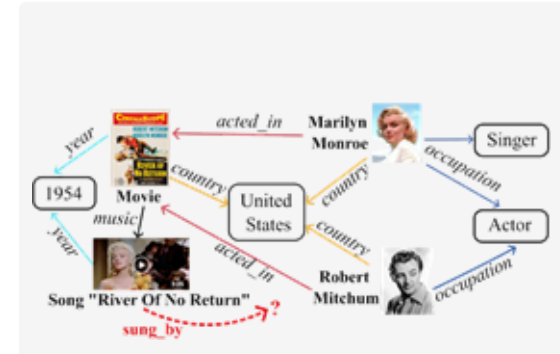
Social Networks



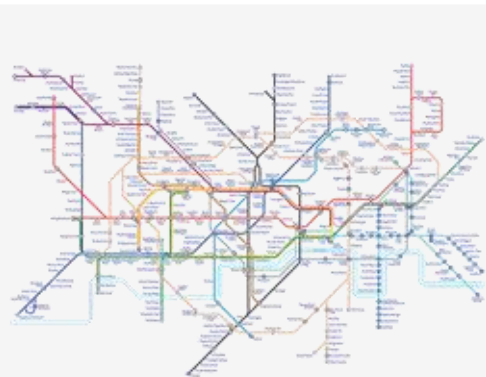
Molecular Networks



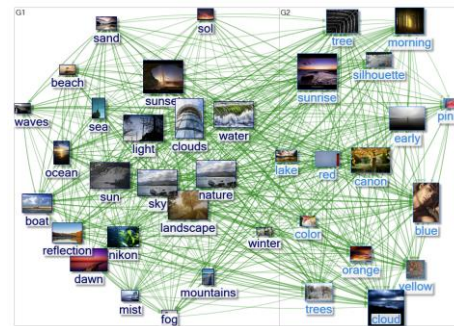
Brain Networks



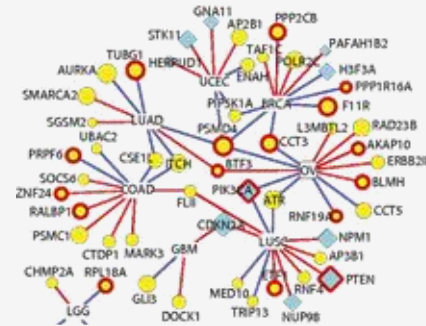
Knowledge Graphs



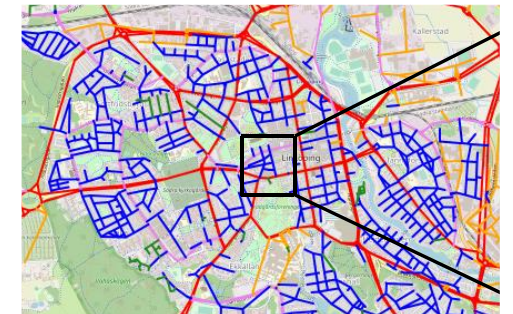
Transportation Networks



Tag Networks



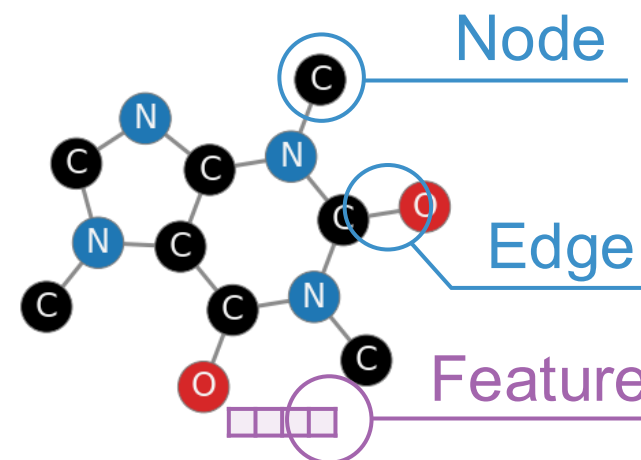
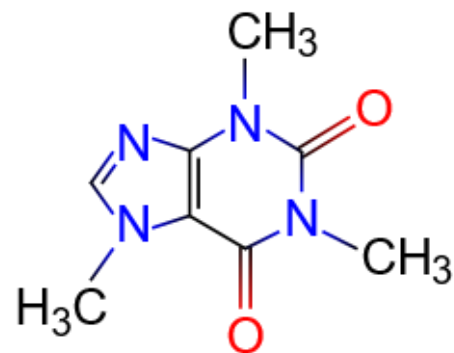
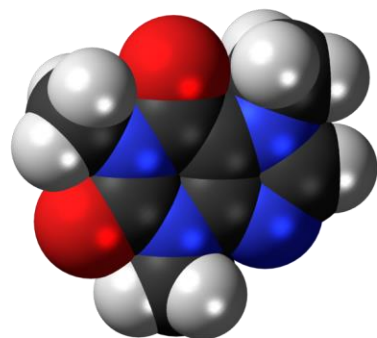
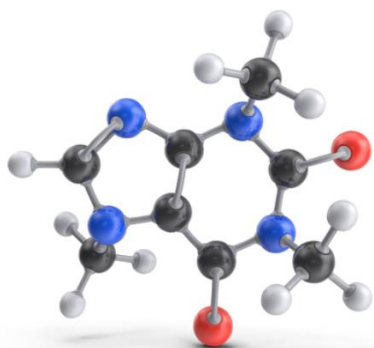
Biological Networks



Road Networks

Graph Representation

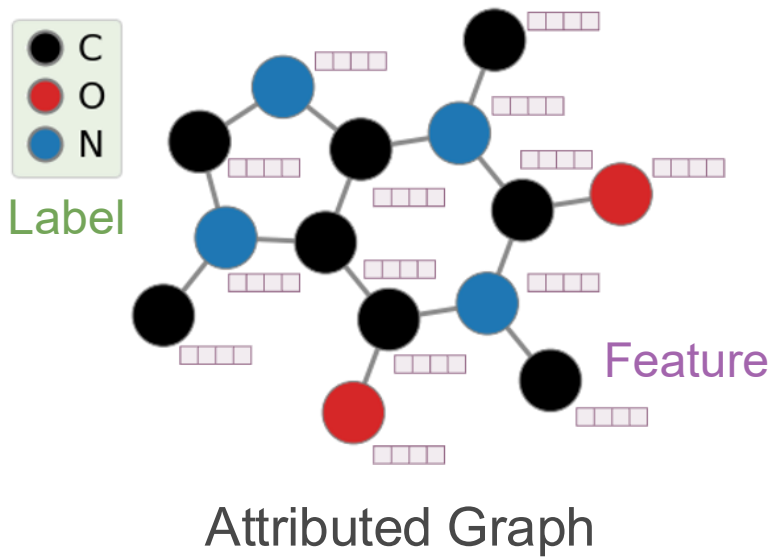
- Graph: entities (nodes) & relationships (edges)
- Attributed Graph: & feature vectors



Different models of the caffeine molecule

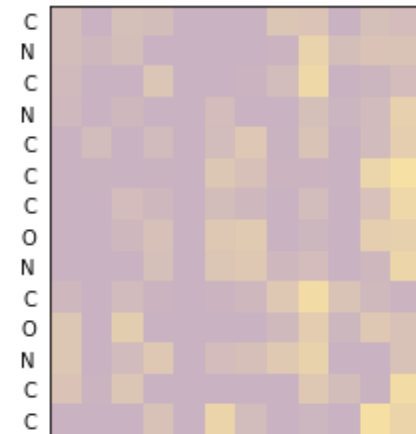
Graph Representation

- Node Level: Feature x , Label y
- Graph Level: Adjacency matrix A , Feature matrix X



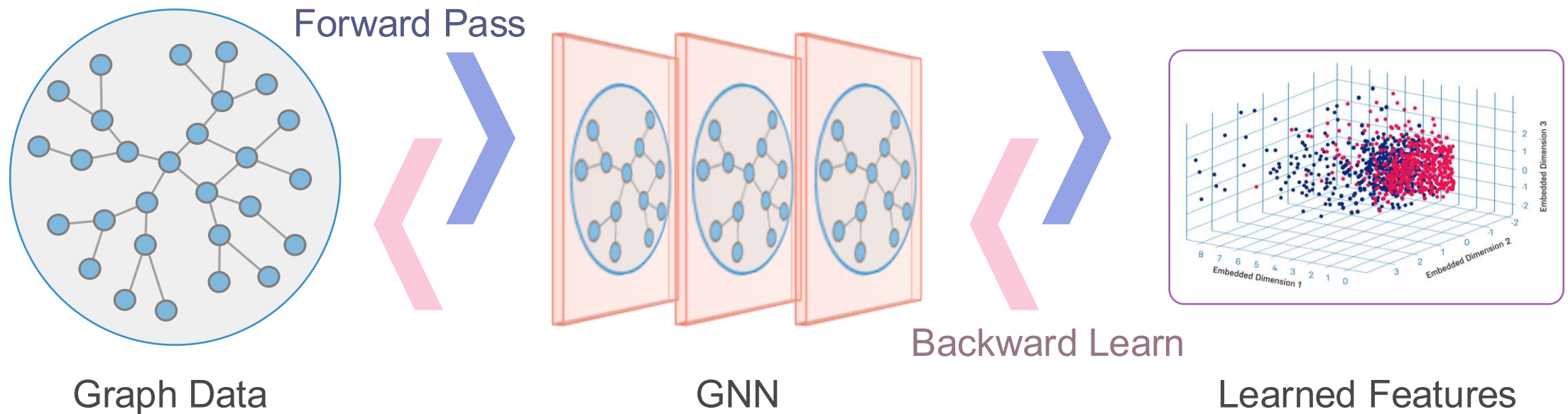
C	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0
C	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0
N	0	0	0	0	1	0	0	0	0	1	0	0	1	0	0
O	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0
C	0	0	0	0	0	0	0	0	0	1	0	1	1	0	0
N	0	0	0	0	0	0	1	0	0	1	0	0	0	0	1
O	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0
C	0	0	0	0	0	1	0	1	1	0	0	0	0	0	0
C	0	1	0	0	1	0	1	0	0	0	0	0	0	0	0
C	0	0	0	1	0	1	0	0	0	0	0	0	1	0	0
N	0	0	1	0	1	0	0	0	0	0	0	0	0	0	0
C	0	1	0	1	0	0	0	0	0	0	0	0	0	0	0
N	1	0	1	0	0	1	0	0	0	0	0	0	0	0	0
C	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0
	C	N	C	N	C	C	C	O	N	C	O	N	C	C	C

Adjacency matrix



GNN: Graph Neural Network

- Design: graph data + neural network
- Data: graph data (irregular) \rightarrow features (structured)



GNN: Graph Neural Network

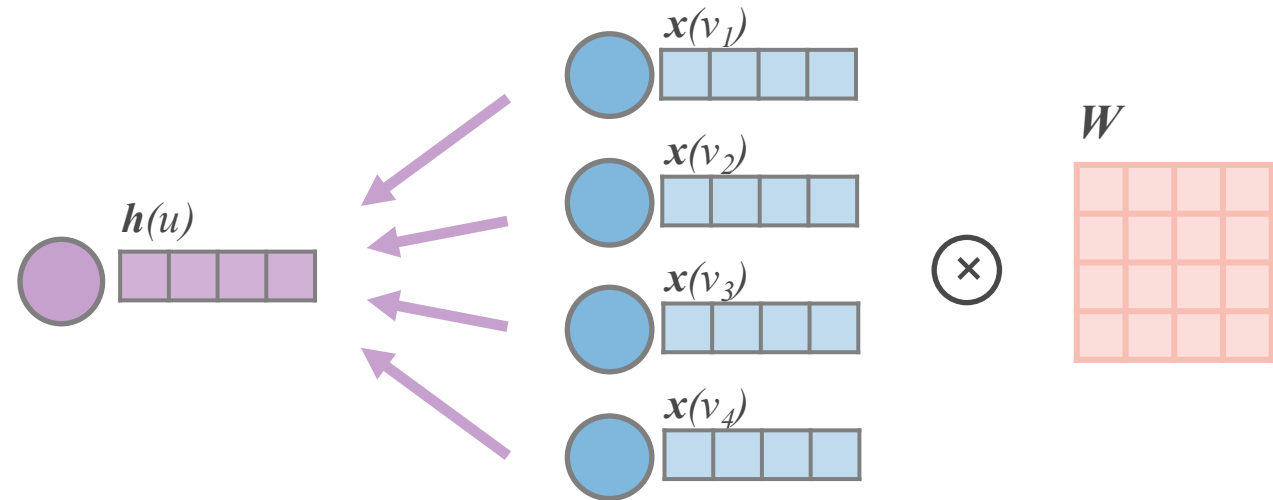
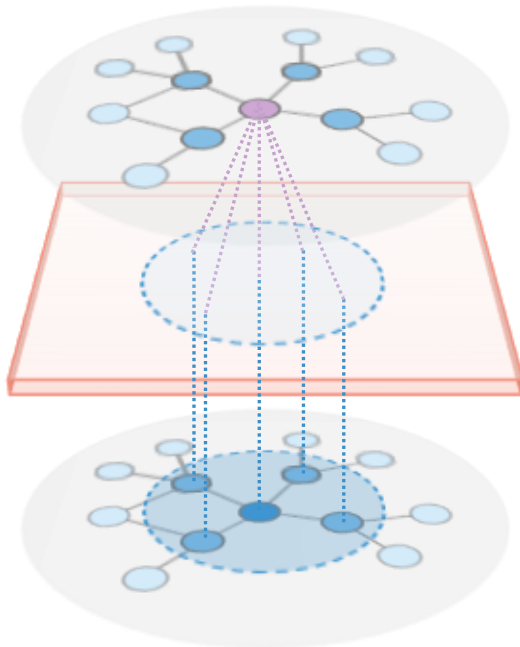
- Graph Convolution:

$$h(u) = \sum_{v \in N(u)} x(v) \cdot W$$

Node Representation

Aggregate Neighbors

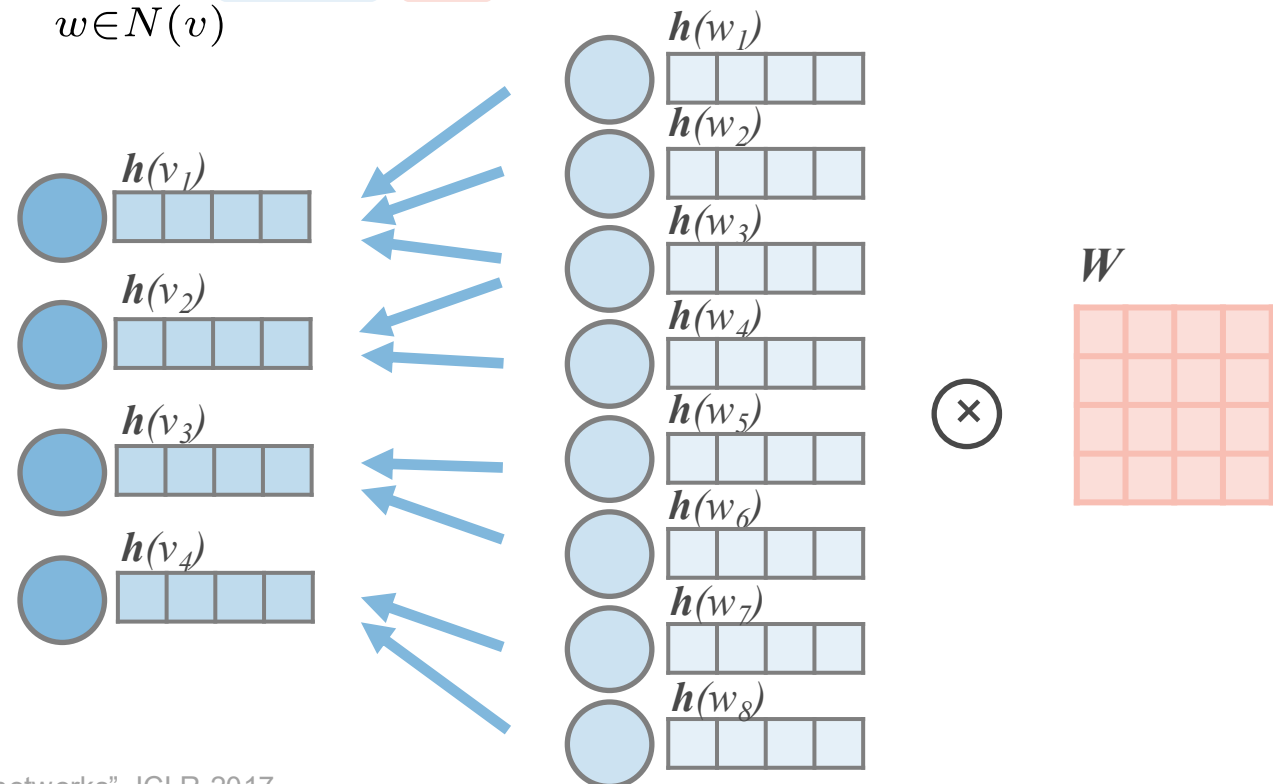
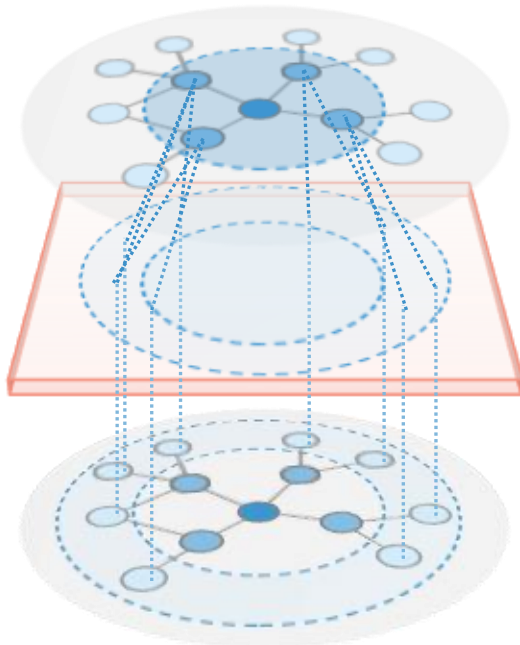
Learnable Weights



GNN: Graph Neural Network

- Stack multi-hop graph convolutions as layers

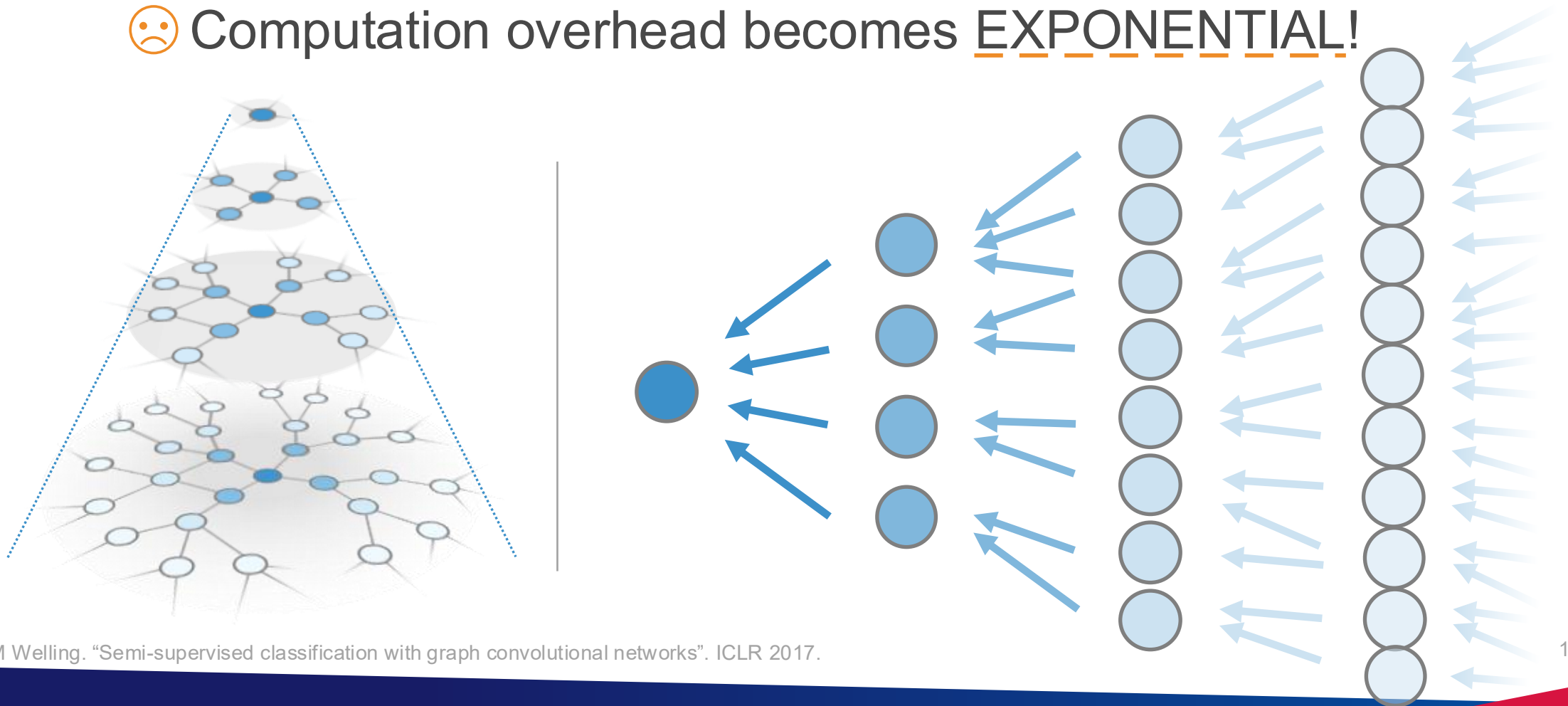
$$\mathbf{h}(v) = \sum_{w \in N(v)} \mathbf{h}(w) \cdot \mathbf{W}$$



Challenge: GNN Scalability Issue

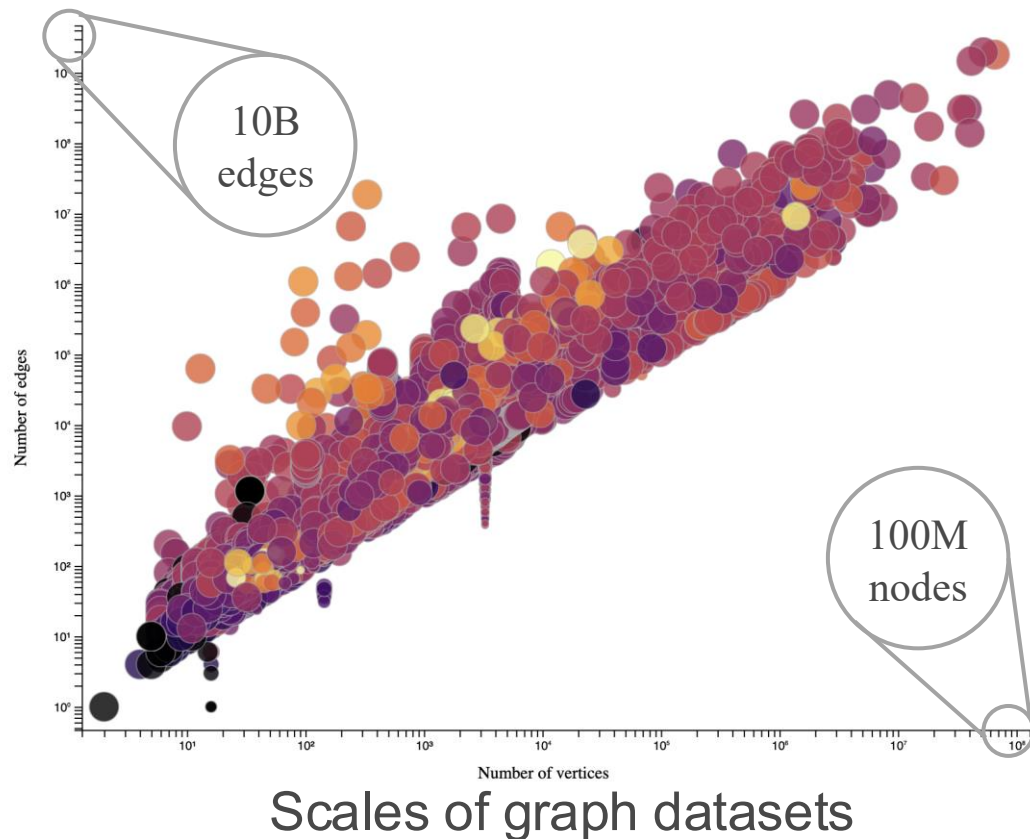
- If keep stacking more layers ...

☹️ Computation overhead becomes EXPONENTIAL!

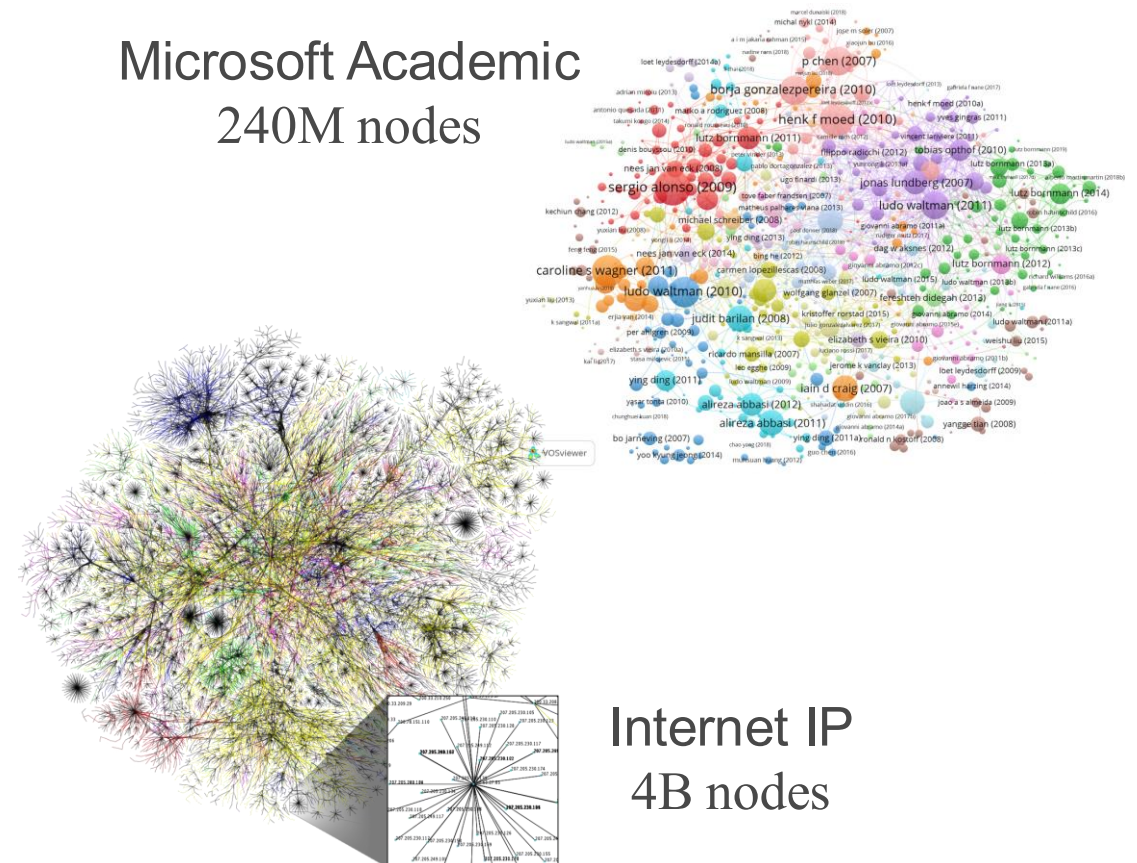


Challenge: GNN Scalability Issue

- Modern real-world graphs are on the scale of millions or billions



Microsoft Academic
240M nodes

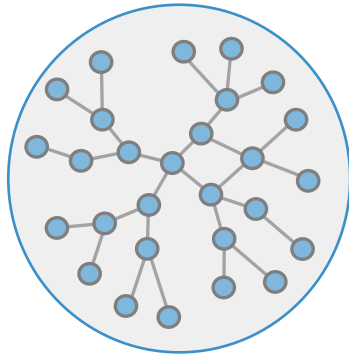


- T Peixoto, "The Netschleuder network catalogue and repository", 2020. <https://networks.skewed.de/>

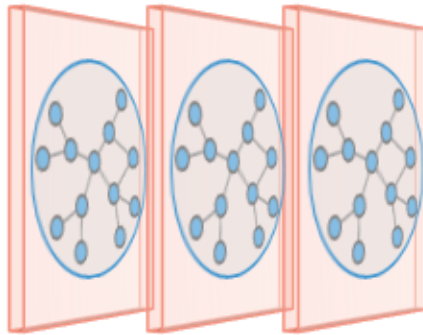
- JM Díaz-Nafría & T Guarda. "Is the Internet-of-Things a Burden or a Leverage for the Human Condition?" MDPI Proceedings. 2017.

Objective: Scaling Up GNN

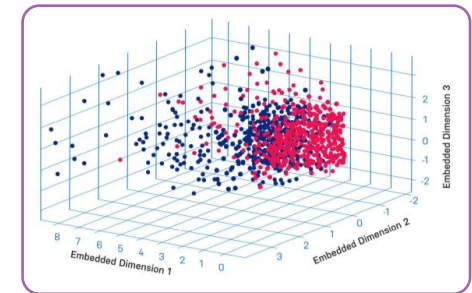
Graph Data



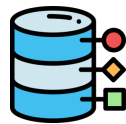
GNN



Learned Features



Deploy Scalable Data Structures



- Transformer model?
- Subgraph model?

Evaluate Performance



- Approximation theory?
- Efficiency benchmark?

Design Scalable Algorithms



- Scalability bottlenecks?
- Long-range convolution?

Outline

Introduction



Data Structures



Evaluations



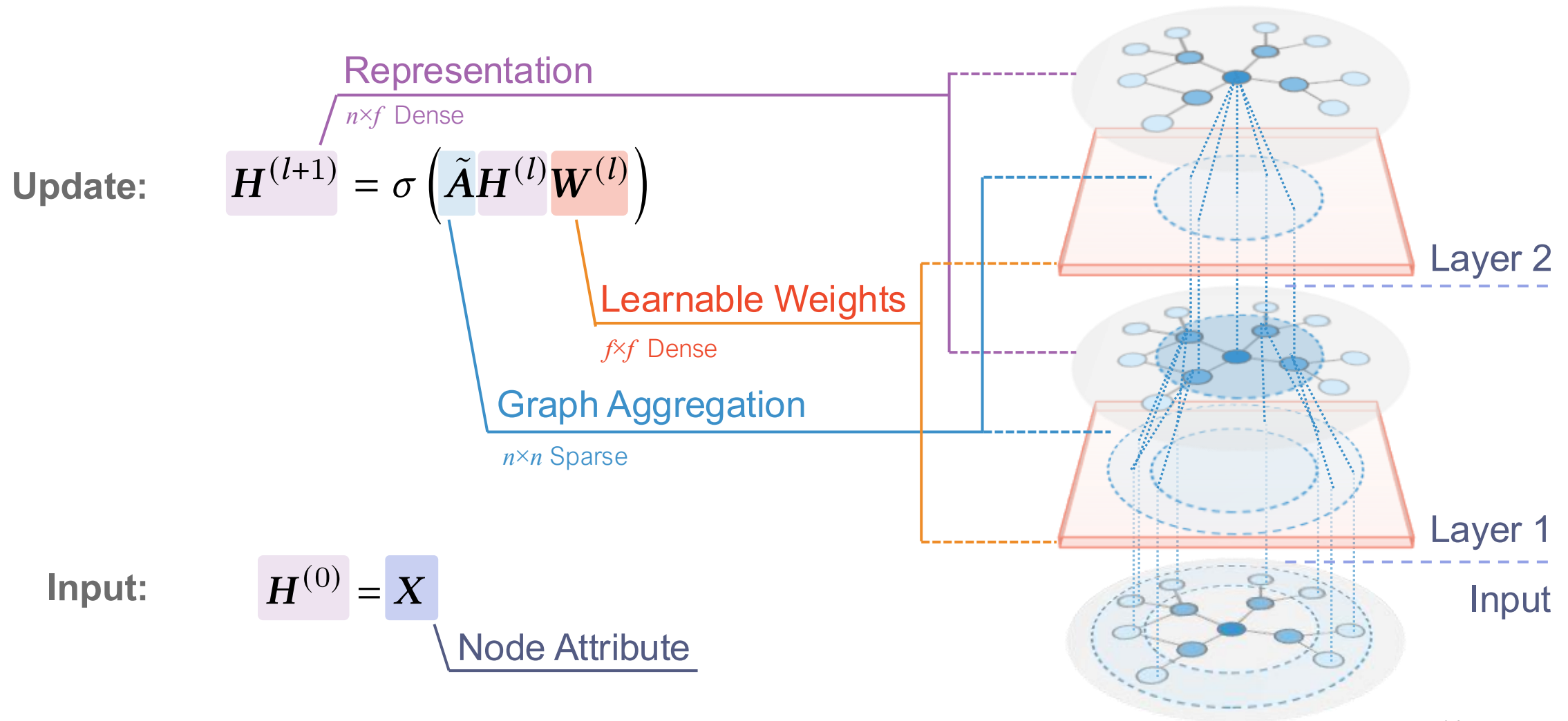
**Design Scalable
Algorithms**

SCARA: Feature-dimension Utilization

LD²: Long Distance & Embedding [Liu et al. 2023]

Conclusion and Future Directions

Background: GCN Complexity



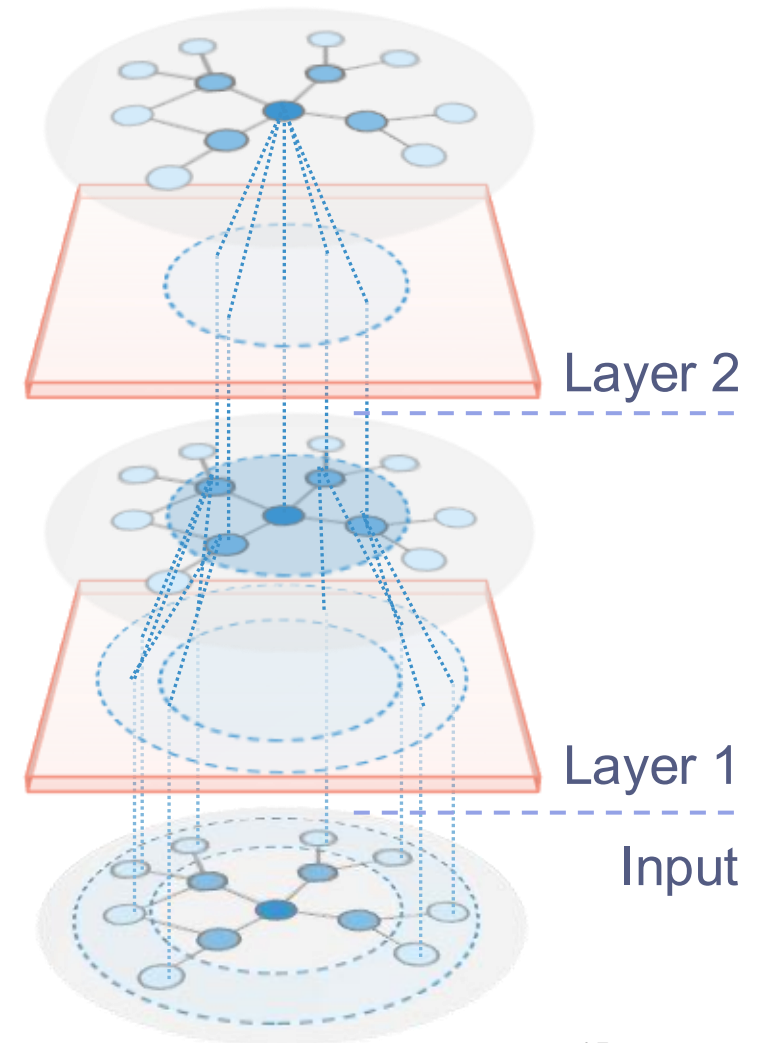
Background: GCN Complexity

Each layer: $H^{(l+1)} = \sigma \left(\tilde{A} H^{(l)} W^{(l)} \right)$

$O(mf)$ Graph Aggregation
Sparse-dense Matrix Multiplication

$O(nf^2)$ Weight Transformation
Dense-dense Matrix Multiplication

Total: $O(Lmf + Lnf^2)$



SCARA: Iterative → Decoupled Computation

- GCN:

Iterative: $L \times$

$$H^{(l+1)} = \sigma \left(\tilde{A} H^{(l)} W^{(l)} \right)$$

- SCARA (ours):

Normalize: $1 \times$

$$\tilde{A}^l \cdot X = D^{r-1} (AD^{-1})^l D^{1-r} X$$

Feature

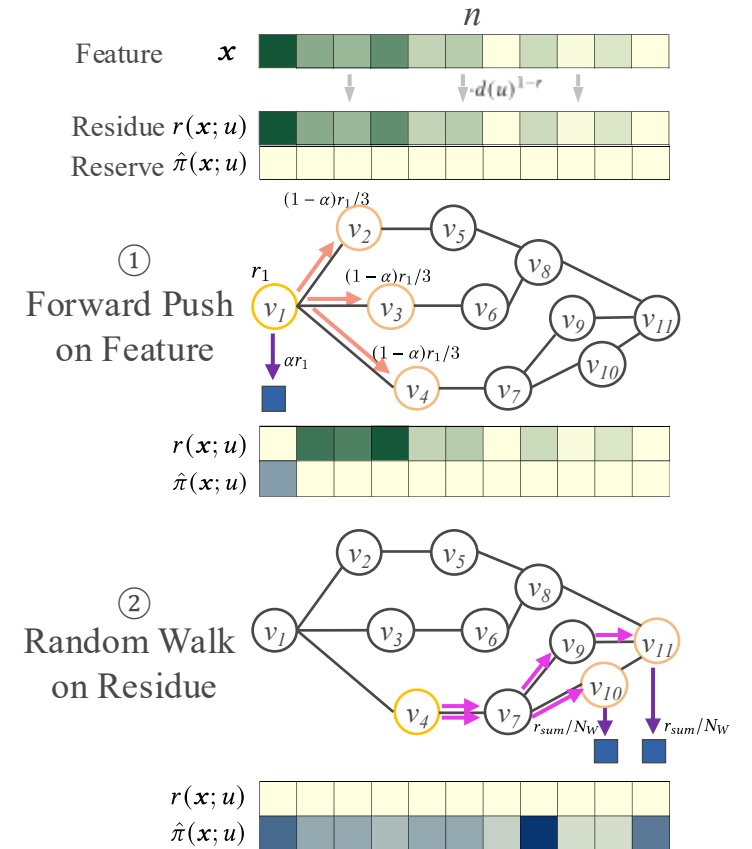
Personalized PageRank

Propagate: $1 \times$

$$H^{(0)} = \sum_0^{\infty} \alpha (1 - \alpha)^l \tilde{A}^l \cdot X$$

Transform: $L \times$

$$H^{(l+1)} = \sigma \left(H^{(l)} W^{(l)} \right)$$



SCARA: Iterative \rightarrow Decoupled Computation

- GCN:

Iterative: $O(LmF + LnF^2)$

☹️ Coupled graph & NN computation

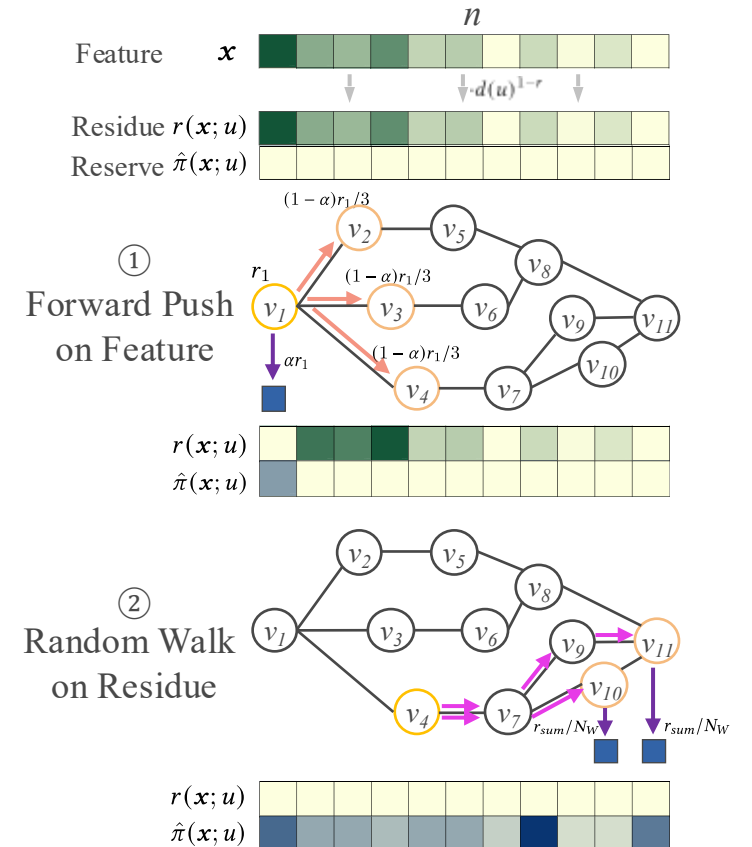
- SCARA (ours):

Propagate: $O(F\sqrt{m \log n}/\lambda)$

😊 Sublinear CPU precomputation

Transform: $O(LnF^2)$

😊 Efficient GPU training



SCARA: Node \rightarrow Feature Dimension

- GCN:

Node-wise:

$$\tilde{\mathbf{A}} \cdot \mathbf{H} = \sum_{u=1}^n \tilde{\mathbf{A}}[u] \cdot \mathbf{H}[u]$$

☹️ *Limited memory locality*

- SCARA:

Decompose:

$$\begin{aligned} \min \text{rank}(\mathbf{X}_B \Theta) + \eta \|\mathbf{Z}\|_1, \\ \text{s.t. } \mathbf{X}_B \Theta + \mathbf{Z} = \mathbf{X}. \end{aligned}$$

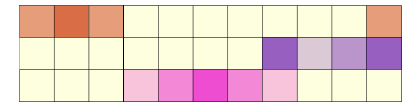
Feature reuse:

$$\check{\mathbf{P}} = \hat{\mathbf{P}}_B \cdot \Theta + \hat{\mathbf{P}}_Z$$

😊 *Reduced compute redundancy*

① Base

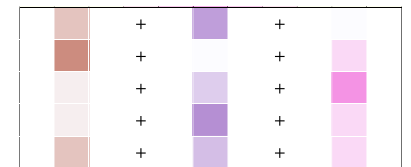
$$\hat{\pi}(\mathbf{b}_i, \gamma \beta_s)$$



Base Calculation
(High Precision)

② Non-base $\pi^*(\mathbf{x}_f)$

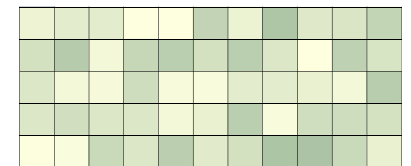
$$\sum \theta_i \cdot \hat{\pi}(\mathbf{b}_i, \gamma \beta_s)$$



Base Combination

+

$$\hat{\pi}(\mathbf{x}', (1 - \gamma \sum \theta_i) \beta_s)$$



Residue Calculation
(Low Precision)

SCARA: Complexity Results

- Time: sub-linear precomputation + efficient learning
- Memory: efficient GPU usage

☹️ *Not suitable for minibatch*

☹️ *Not scalable to large graph m and n*

☹️ *Insufficient memory*

Model	Precomp. Time	Training Time	Precomp. Mem.	Training Mem.
GCN [3]	–	$O(ILmF + ILnF^2)$	–	$O(LnF + LF^2)$
Cluster-GCN [22]	$O(m)$	$O(ILmF + ILnF^2)$	$O(n)$	$O(Ln_bF + LF^2)$
GraphSAINT [23]	–	$O(IL_PLnF^2)$	–	$O(L_PLn_bF + LF^2)$
GAS [24]	$O(m + LnF)$	$O(ILmF + ILnF^2)$	$O(LnF)$	$O(Ldn_bF + LF^2)$
APPNP [25]	$O(m)$	$O(IL_PmF + ILnF^2)$	$O(m)$	$O(Ln_bF + LF^2 + nn_b)$
PPRGo [26]	$O(m/\delta)$	$O(IKnF + ILnF^2)$	$O(n/\delta)$	$O(Ln_bF + LF^2 + Kn_b)$
SGC [27]	$O(L_PmF)$	$O(ILnF^2)$	$O(m)$	$O(Ln_bF + LF^2)$
GBP [11]	$O(L_PF\sqrt{L_Pm\log(L_Pn)}/\epsilon)$	$O(ILnF^2)$	$O(nF)$	$O(Ln_bF + LF^2)$
SCARA (ours)	$O(F\sqrt{m\log n}/\lambda)$	$O(ILnF^2)$	$O(nF)$	$O(Ln_bF + LF^2)$

😊 *Sub-linear precomputation*

😊 *Decoupled simple batching*

SCARA: Experimental Results

- Time: 30-800× faster parallel precomputation
- Memory: Paper100M with 72GB without OOM
- Effectiveness: similar or better F1-score, fast convergence

	Dataset	Nodes n	Edges m	Features F		Dataset	Nodes n	Edges m	Features F		
	Reddit [15]	232,965	114,615,892	602		Papers100M [16]	111,059,956	1,615,685,872	128		
Transductive	Reddit					Papers100M					
	Learn	(Pre. + Train)	Infer	Mem.	F1	Learn	(Pre. + Train)	Infer	Mem.	F1	
GraphSAINT	14.4	(- 14.4)	166.2	13.7	41.6 ± 4.8	-	-	-	-	OOM	-
GAS	1151	(- 1151)	2.2	14.0	38.2 ± 0.3	-	-	-	-	OOM	-
PPRGo	79.4	(62.3 + 17.1)	29.1	9.4	41.5 ± 2.3	-	-	-	-	OOM	-
GBP	138	(124 + 13.7)	13.5	7.9	38.8 ± 0.3	-	-	-	-	OOM	-
SCARA (ours)	13.9	(0.07 + 13.8)	10.7	5.6	44.1 ± 0.4	1346	(12.7 + 1333)	4.7	71.4	35.7 ± 0.9	

Outline

Introduction



Data Structures



Evaluations



**Design Scalable
Algorithms**

SCARA: Feature-dimension Utilization

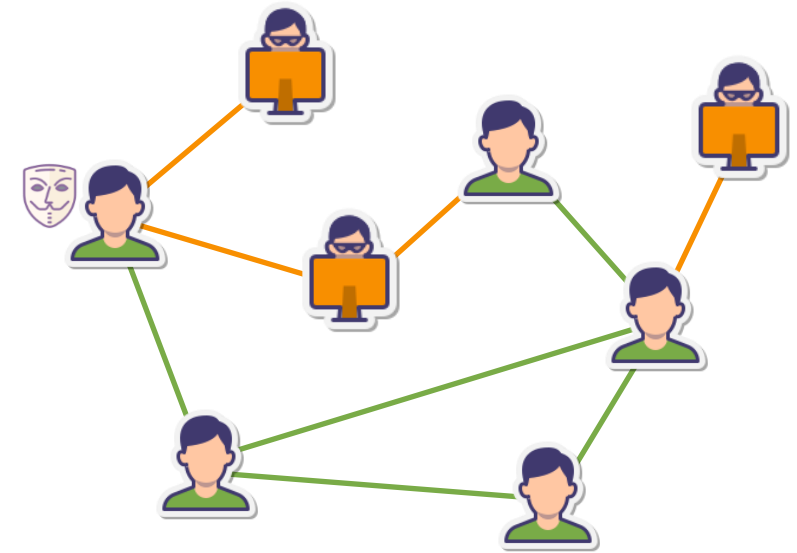
LD²: Long-distance Embedding

[NeurIPS 2023]

Conclusion and Future Directions

Background: Graph Heterophily

- Heterophily: connected nodes tend to be of dissimilar classes
- Graph: financial transaction network with **benign** or **fraud** users
- Task: identify fraud users
- **Fraudsters** use **benign** accounts as springboards 🎭 to reach other **fraudsters**



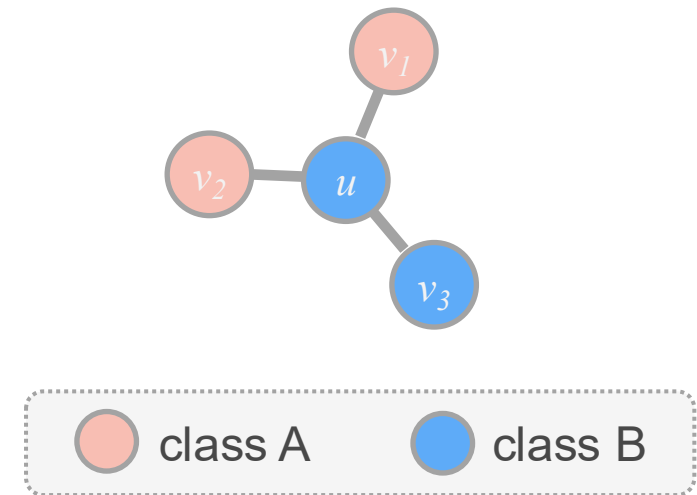
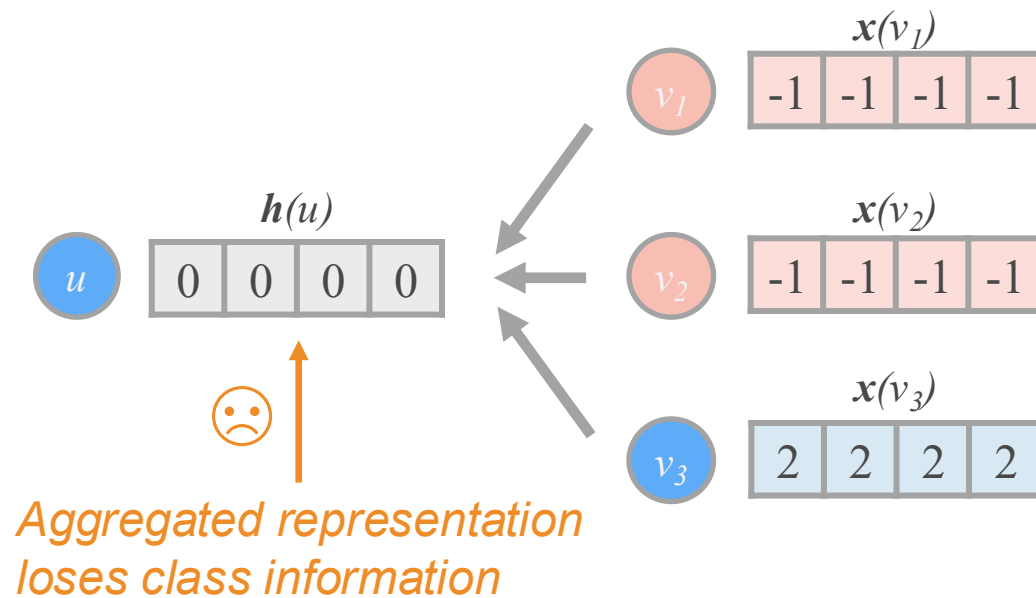
☹️ *1-hop connection is not enough!*



Background: Graph Heterophily

- Heterophily: 1-hop graph convolution fails

How to design decoupled representation under heterophily?



LD²: Representing Heterophily

- SGC:

$$P = \tilde{A}^L \cdot X$$

Local Aggregation

- LD² (ours):

2-hop Adjacency:

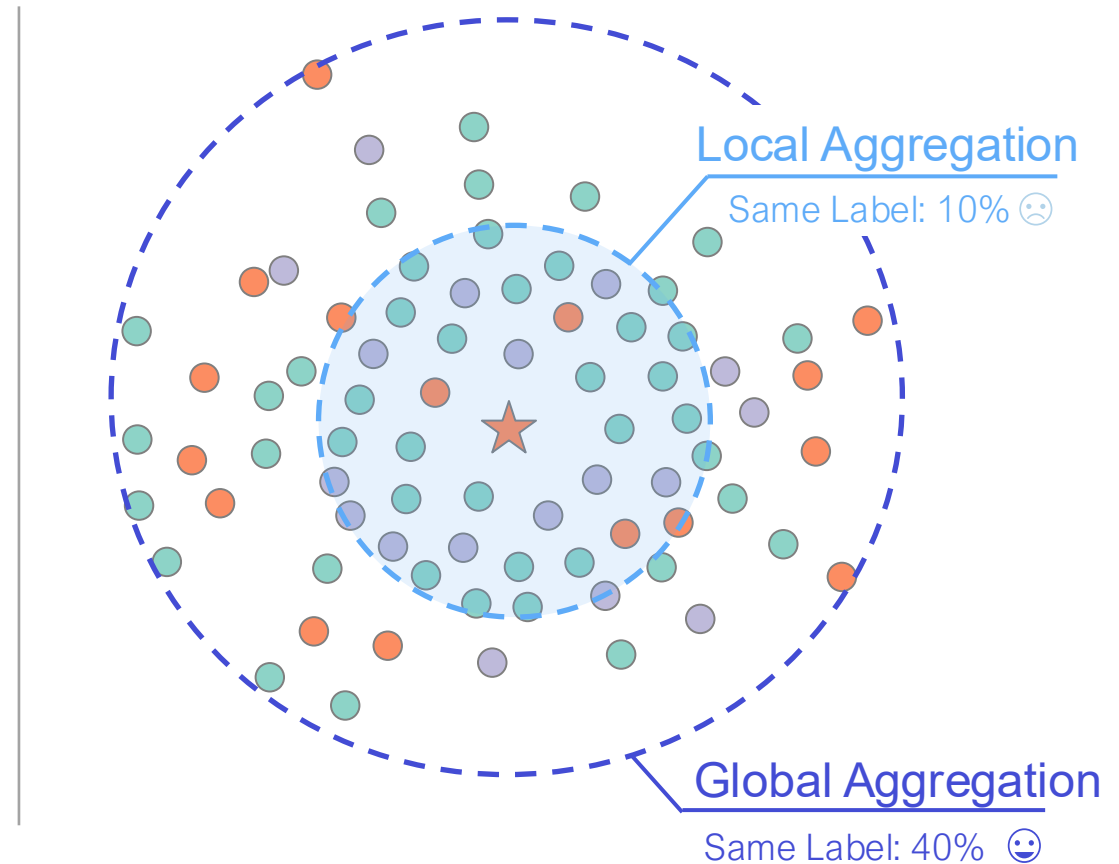
$$P_A = \phi(A^{2L} \cdot \mathcal{N})$$

Constant Feature:

$$P_{X,L2} = \sum_{l=1}^L \bar{A}^{2l} \cdot X$$

Inverse Feature:

$$P_{X,H} = \sum_{l=1}^L \tilde{L}^l \cdot X$$



LD²: Representing Heterophily

- LD² (ours):

2-hop Adjacency:

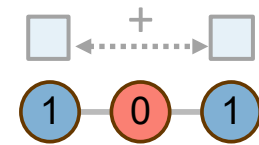
$$P_A = \phi(A^{2L} \cdot \mathcal{N})$$

 Spatial

$$P_A \cdot P_A^T \approx A^2$$

Low-rank approximation

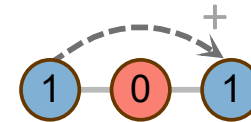
 Spectral



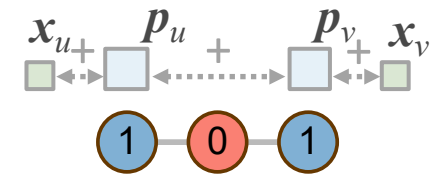
Structure smoothing

Constant Feature:

$$P_{X,L2} = \sum_{l=1}^L \bar{A}^{2l} \cdot X$$



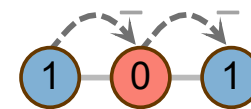
Long-range similarity



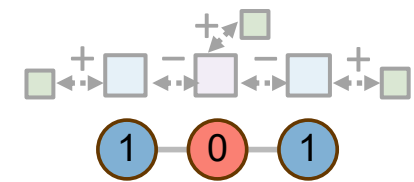
Feature smoothing

Inverse Feature:

$$P_{X,H} = \sum_{l=1}^L \tilde{L}^l \cdot X$$



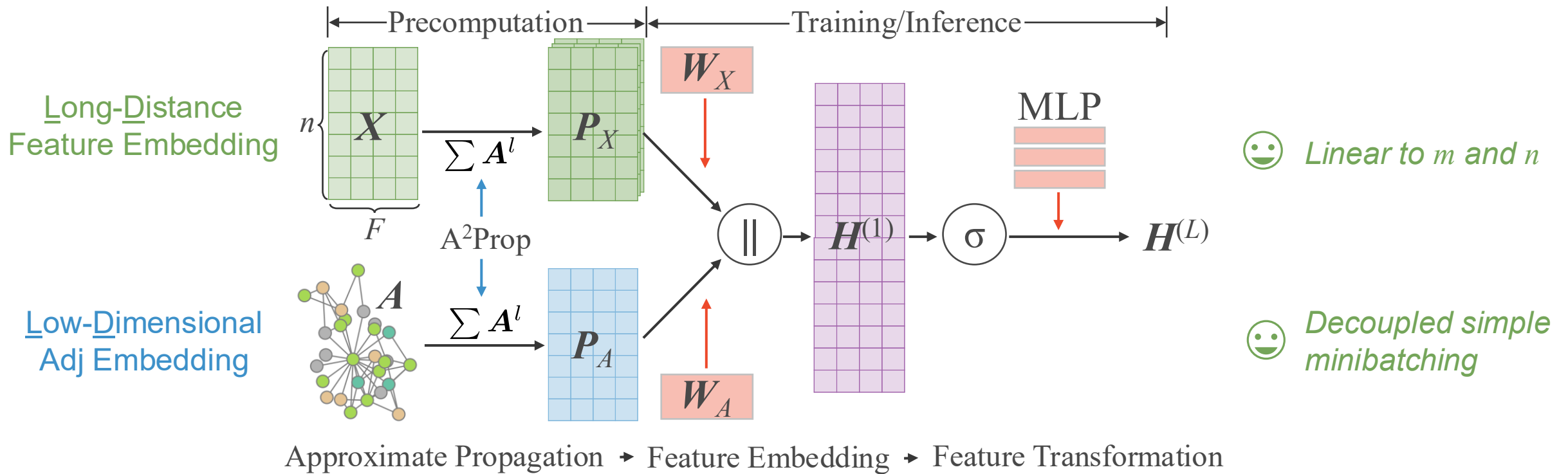
Fine-grained details



Structure difference

LD²: Decoupled Architecture

- Spectral precomputation + simple feature transformation



LD²: Complexity Results

- Minibatching without graph-scale overhead

Model	Precomp. Time	Training Time	Inference Time
GPRGNN [39]	$O(m)$	$O(IL_P m F + ILnF^2)$	$O(L_P m F + LnF^2)$
GCNJK [40]	–	$O(ILmF + ILnF^2)$	$O(LmF + LnF^2)$
MixHop [41]	–	$O(IL_P LmF + ILnF^2)$	$O(L_P LmF + LnF^2)$
LINKX [42]	–	$O(ImF + ILnF^2)$	$O(mF + LnF^2)$
LD² (ours)	$O(L_P m F)$	$O(ILnF^2)$	$O(LnF^2)$

☹️ *Not scalable to large graph m and n*

😊 *Linear to m and n*

Model	Precomp. Mem.	Training Mem.	Inference Mem.
GPRGNN [39]	$O(m)$	$O(LnF + LF^2 + m)$	$O(LnF + LF^2 + m)$
GCNJK [40]	–	$O(L_C n F + L_C F^2)$	$O(L_C n F + L_C F^2)$
MixHop [41]	–	$O(CLnF + CLF^2)$	$O(CLnF + CLF^2)$
LINKX [42]	–	$O(L_C n_b F + L_C F^2 + nF)$	$O(L_C n_b F + L_C F^2 + nF)$
LD² (ours)	$O(CnF)$	$O(L_C n_b F + L_C F^2)$	$O(L_C n_b F + L_C F^2)$

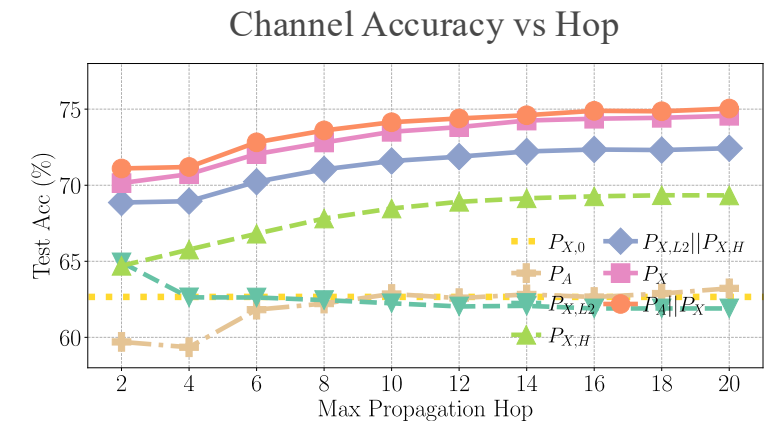
☹️ *Not suitable for minibatch*

😊 *Decoupled simple batching*

LD²: Experimental Results

- Efficiency: 3-15× faster, 5× lower memory
- Effectiveness: No accuracy drop for minibatch.

Dataset	twitch-gamers			pokec			snap-patents			wiki		
	Learn	Infer	Mem.	Learn	Infer	Mem.	Learn	Infer	Mem.	Learn	Infer	Mem.
MLP	6.36	0.02	0.61	47.86	0.11	13.77	27.39	0.28	9.33	133.55	0.62	18.15
PPRGo	10.46+15.88	0.41	9.64	121.95+56.11	2.69	3.82	(>12h)			(>12h)		
SGC	0.09+0.74	0.01	0.28	1.05+8.08	0.01	0.28	4.94+23.54	0.01	0.42	12.66+7.98	0.01	0.52
GCNJK-GS	71.48	0.02*	7.33	27.33	0.09*	9.03	19.02	0.23*	9.21	95.52	0.69*	16.36
MixHop-GS	52.12	0.01*	1.49	71.35	0.03*	12.91	45.24	0.16*	19.58	84.22	0.23*	16.28
LINKX	10.99	0.19	2.35	28.77	0.33	9.03	39.80	0.22	21.53	180.71	1.14	14.53
LD² (ours)	0.85+1.96	0.01	1.44	17.95+6.18	0.01	3.82	31.32+6.96	0.02	3.96	28.12+6.50	0.01	4.47



Takeaway

SCARA: Feature-dimension Utilization

- Design: Decouple graph convolution as PPR
- Algorithm: Reuse features during graph computation
- Experiment: 20 mins for billion-edge graphs

LD²: Long-distance Embedding

- Algorithm: Adapt decoupled propagation under heterophily
- Theory: Long-range spectral aggregation
- Experiment: Minibatching on large-scale heterophilous graphs

Outline

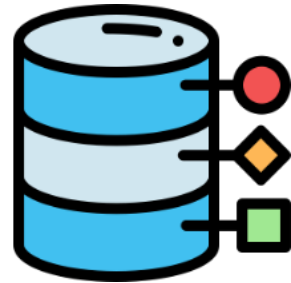
Introduction



Algorithms



Evaluations



**Deploy Scalable
Data Structures**

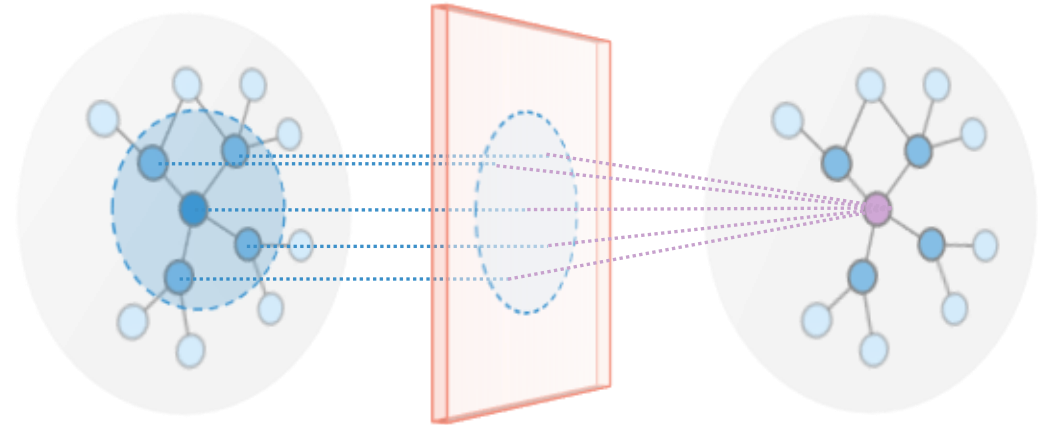
HubGT: Hierarchical Hub Labeling
GENTILE: Subgraph Streaming
[NeurIPS 2023]

Conclusion and Future Directions

Background: Subgraph GNN

- Full-graph learning:

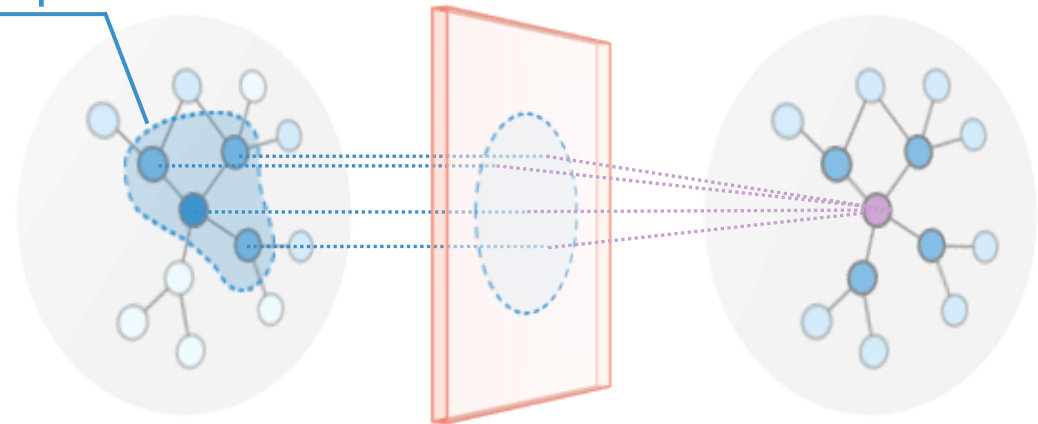
$$\mathbf{h}(u) = \sum_{v \in N(u)} \mathbf{x}(v) \cdot \mathbf{W}$$



- Subgraph learning:

$$\mathbf{h}(u) = \sum_{v \in S(u)} \mathbf{x}(v) \cdot \mathbf{W}$$

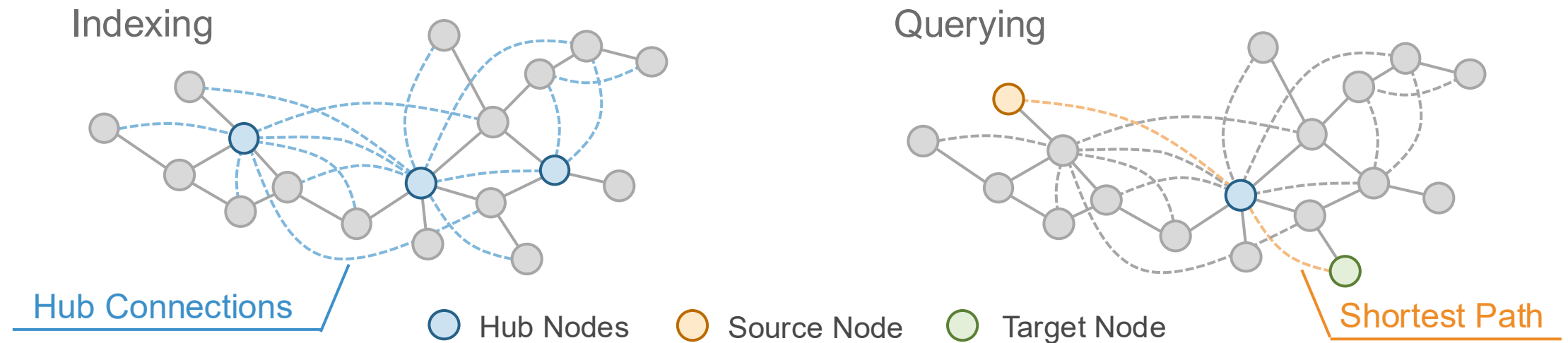
Subgraph



😊 *Smaller graph data*

HubGT: Hub Labeling

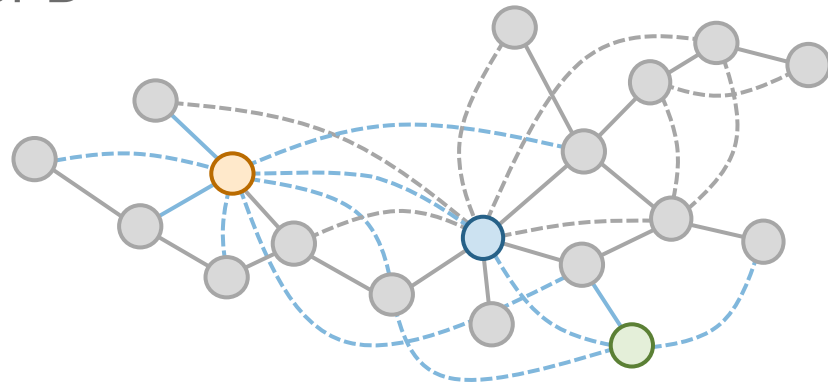
- Hub Labeling: adding links to selected **hub** nodes → distant node pairs can be easily reached through **hubs**
- **Shortest Path Distance (SPD)** can be computed through **hubs**



HubGT: Hub Labeling

- Subgraph SPD (ours): querying any node pairs within a subgraph

Full-graph
SPD



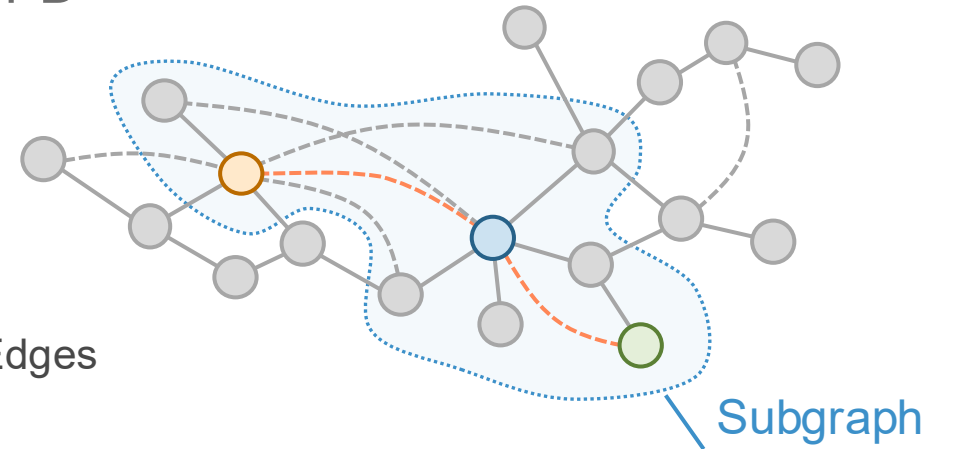
Query Time

$O(s)$

Index Size

$O(ns^2)$

Subgraph
SPD

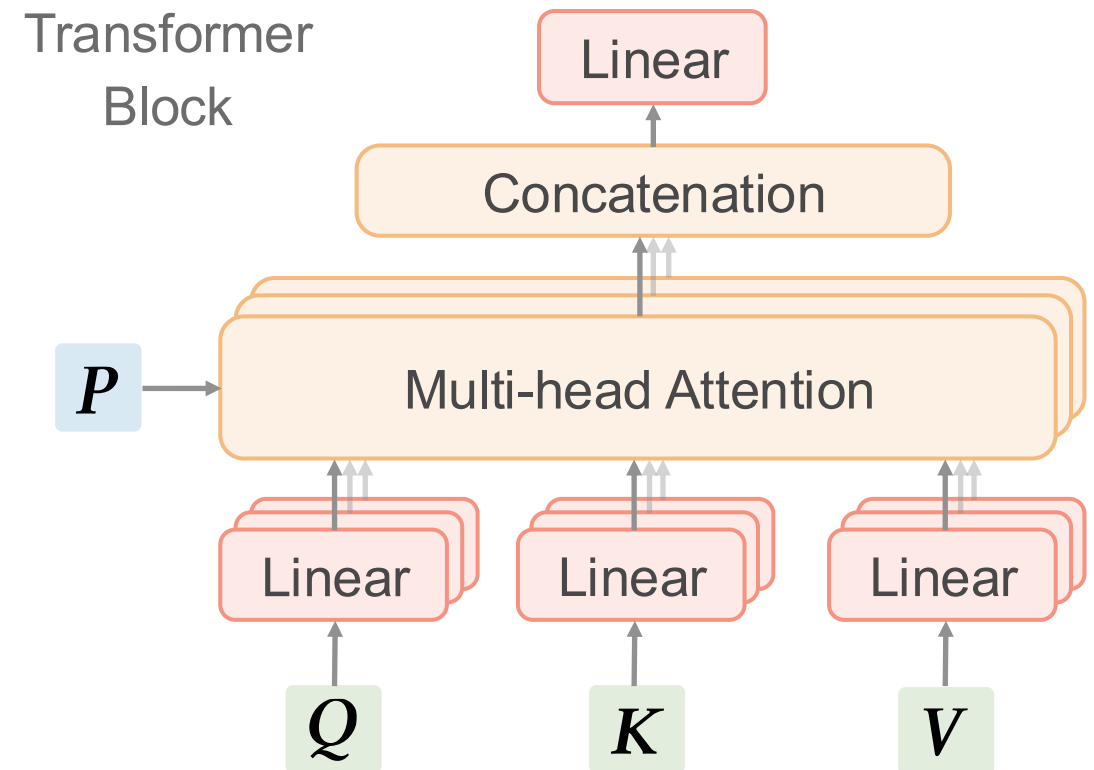
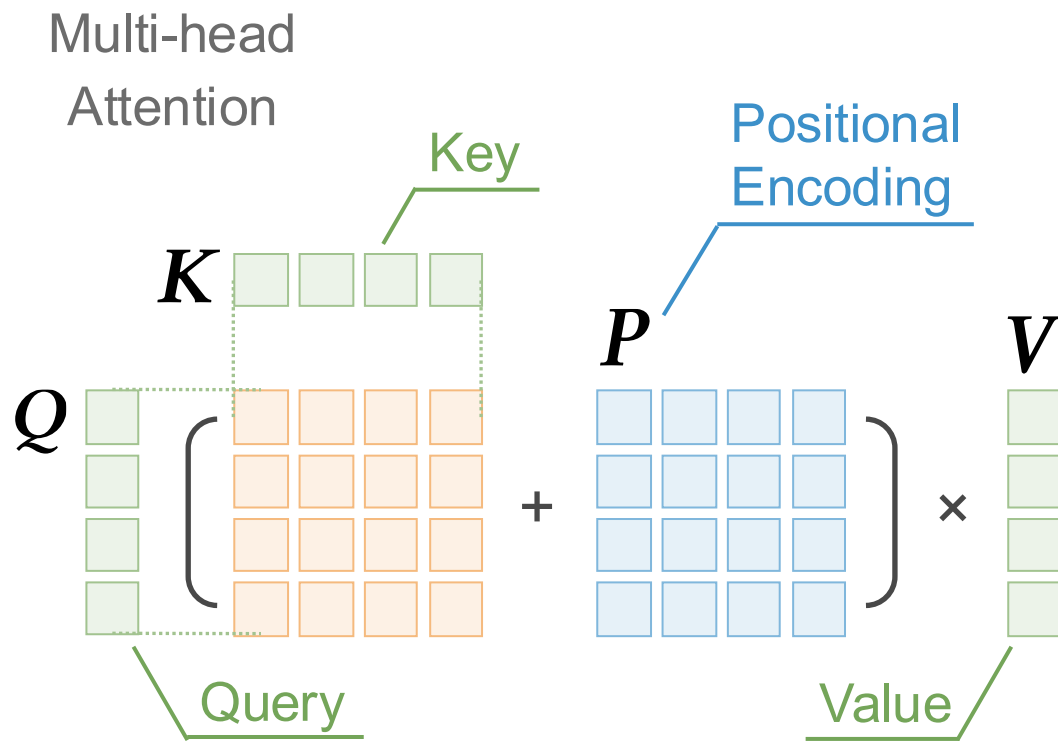


$O(1)$

$O(ns^2)$

Background: Transformer PE

- Positional Encoding (PE): relation as bias in attention



HubGT: Graph Transformer PE

- Positional Encoding: emphasize important pair-wise relation

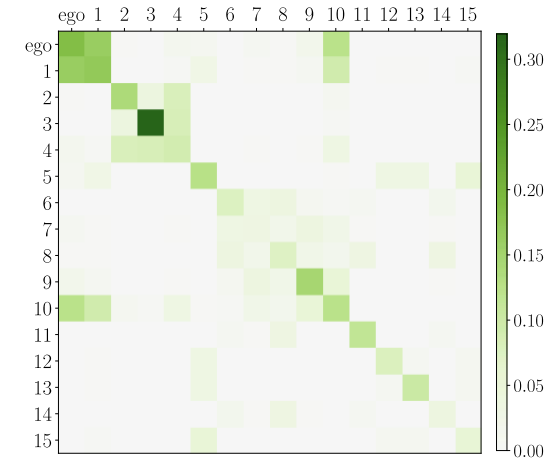
$$\tilde{H}_i = \text{softmax} \left(\frac{Q_i K_i^\top}{\sqrt{F_K}} + P \right) V_i,$$

Node Feature \tilde{H}_i
Positional Encoding P
Transformer Attention $\frac{Q_i K_i^\top}{\sqrt{F_K}}$
Transformer Attention V_i

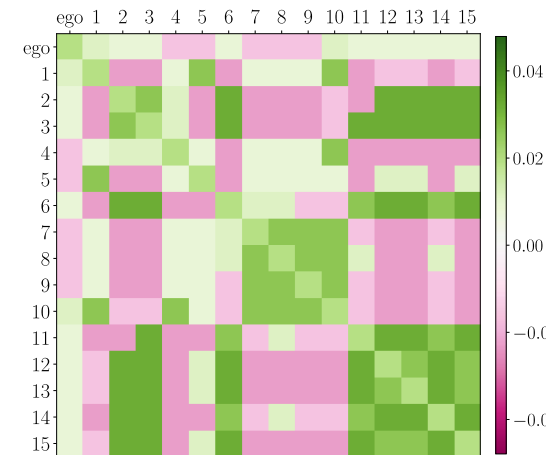
- Complexity:

$O(n^2)$
Full graph

$O(s^2)$
Subgraph



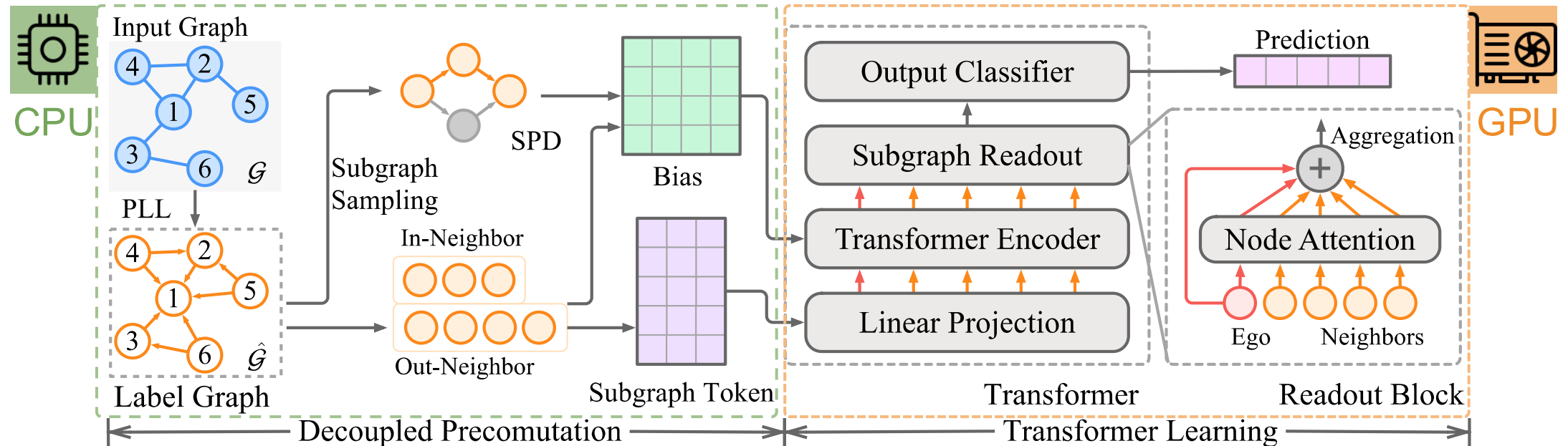
Adjacency
PE



SPD
Encoding

HubGT: Graph Transformer PE

- Hub node \rightarrow hierarchical subgraph \rightarrow input token
- SPD value \rightarrow positional encoding \rightarrow attention bias



HubGT: Complexity Results

- Only linear time & memory complexity

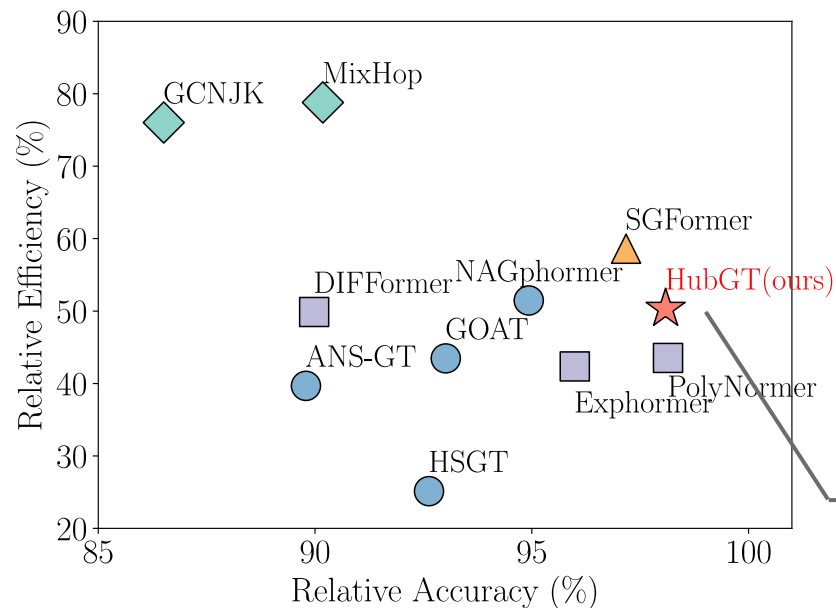
☹️ *Not scalable to large m and n*

☹️ *Insufficient memory*

Taxonomy	Model	Precompute Time	Train Time	RAM Mem.	GPU Mem.	Scale
Vanilla (FB)	Graphormer [2]	$O(n^3)$	$O(Ln^2F)$	$O(n^2)$	$O(Ln^2F)$	0.3K/0.6K
	GRPE [18]	$O(n^2)$	$O(Ln^2F)$	$O(n^2)$	$O(Ln^2F)$	0.3K/0.6K
Kernel-based (NS)	GraphGPS [7]	$O(n^3)$	$O(LnF^2 + LmF)$	$O(nF + n^2)$	$O(Ln_bF + m)$	1.0K/3.0K
	Exphormer [19]	$O(dn)$	$O(LnF^2 + LmF)$	$O(nF + m)$	$O(Ln_bF + m)$	0.2M/1.2M
	NodeFormer [6]	–	$O(LnF^2 + LmF)$	$O(nF + m)$	$O(Ln_bF + m)$	2.4M/60M
	DIFFormer [8]	–	$O(LnF^2 + LmF)$	$O(nF + m)$	$O(Ln_bF + m)$	1.6M/40M
	PolyNormer [12]	–	$O(LnF^2 + LmF)$	$O(nF + m)$	$O(Ln_bF + m)$	2.4M/0.1B
Hierarchical (RS)	NAGphormer [9]	$O(LmF_0)$	$O(LnF^2)$	$O(LnF)$	$O(Ln_bF^2)$	2.4M/60M
	PolyFormer [11]	$O(LmF_0)$	$O(LnF^2)$	$O(LnF)$	$O(Ln_bF^2)$	0.2M/30M
	ANS-GT [13]	$O(ns^2 + md^L)$	$O(LnF^2 + Lns^2F + Lm)$	$O(nF + ns^2 + md^L)$	$O(Ln_bF + n_b s^2)$	20K/0.2M
	GOAT [10]	$O(nF)$	$O(LnF^2 + LmF)$	$O(nF + m)$	$O(Ln_b^2 + Ln_bF + m)$	2.9M/0.1B
	HSGT [14]	$O(n + md^L)$	$O(LnF^2 + LmF)$	$O(nF + md^L)$	$O(Ln_b^2 + Ln_bF + Lm)$	2.4M/0.1B
	HubGT (ours)	$O(ns^3 + ms)$	$O(LnF^2 + Lns^2F)$	$O(nF + ns^2)$	$O(Ln_bF + n_b s^2)$	1.6M/0.1B

HubGT: Experimental Results

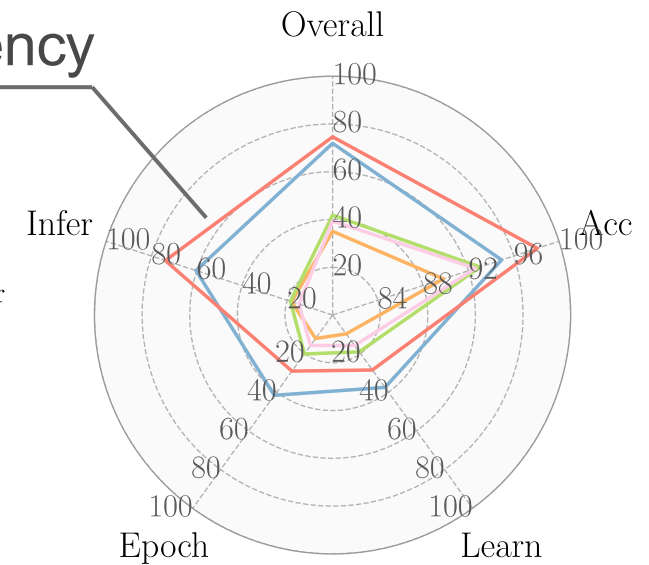
- Efficiency: fast & scalable minibatch
- Effectiveness: higher accuracy under heterophily



Fast + top-tier Accuracy

Overall Efficiency

- Expformer
- DIFFormer
- PolyNormer
- NAGphormer
- ANS-GT
- GOAT
- HSGT
- HubGT



Outline

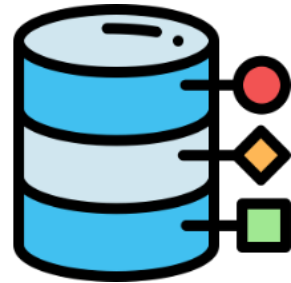
Introduction



Algorithms



Evaluations



**Deploy Scalable
Data Structures**

HubGT: Hierarchical Hub Labeling

GENTI: Subgraph Streaming

[VLDB 2024]

Conclusion and Future Directions

Background: Walk-based GNN

- Subgraph by convolution:

$$P[u, v] = H[\mathcal{N}(u)] + H[\mathcal{N}(v)]$$

Link
Feature

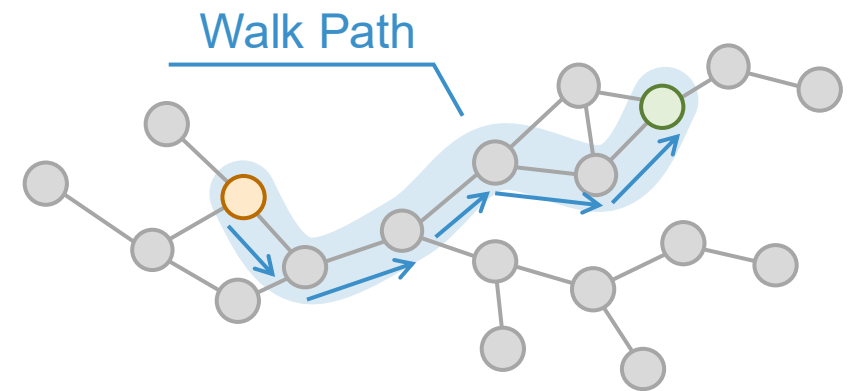
☹️ *Include redundant nodes*



- Subgraph by random walks:

$$P[u, v] = H[\{u, t_1, t_2, \dots, v\}]$$

😊 *More relevant nodes*

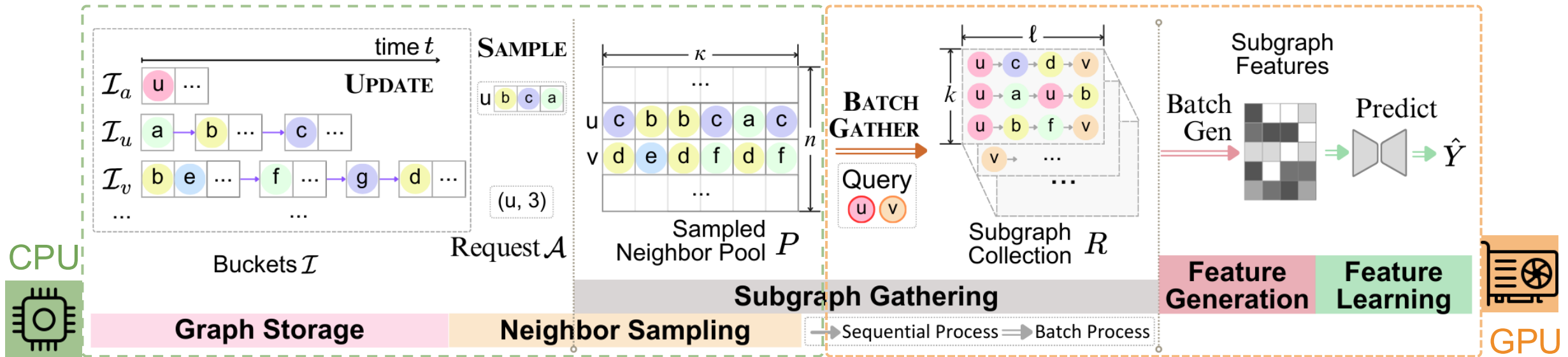


○ Source Node ○ Target Node

○ Insignificant Nodes

GENTI: Subgraph GNN Operators

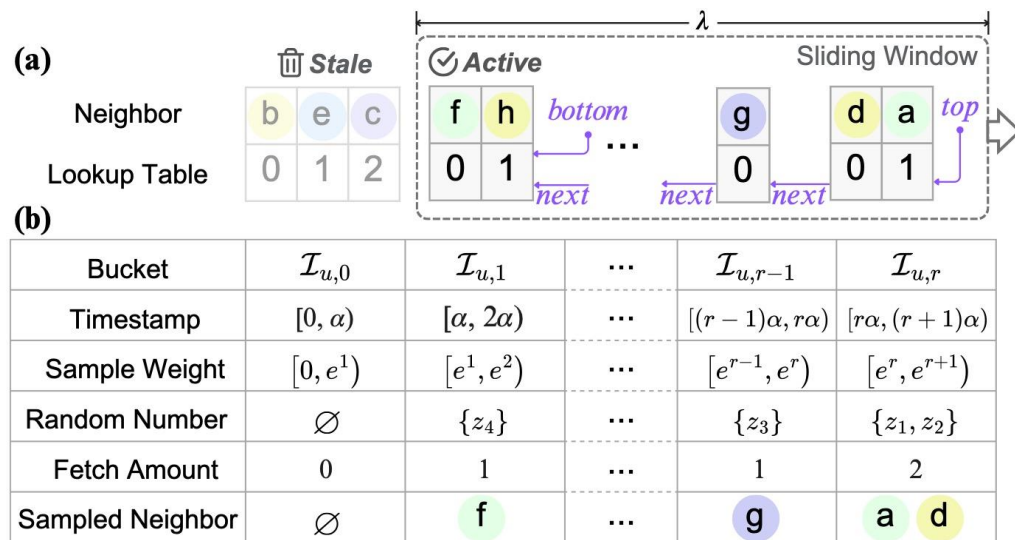
- Sampling: linked buckets arranged in temporal order
- Gathering: pipelining sample (CPU) and gather (GPU)



GENTI: Streaming in Subgraph GNN

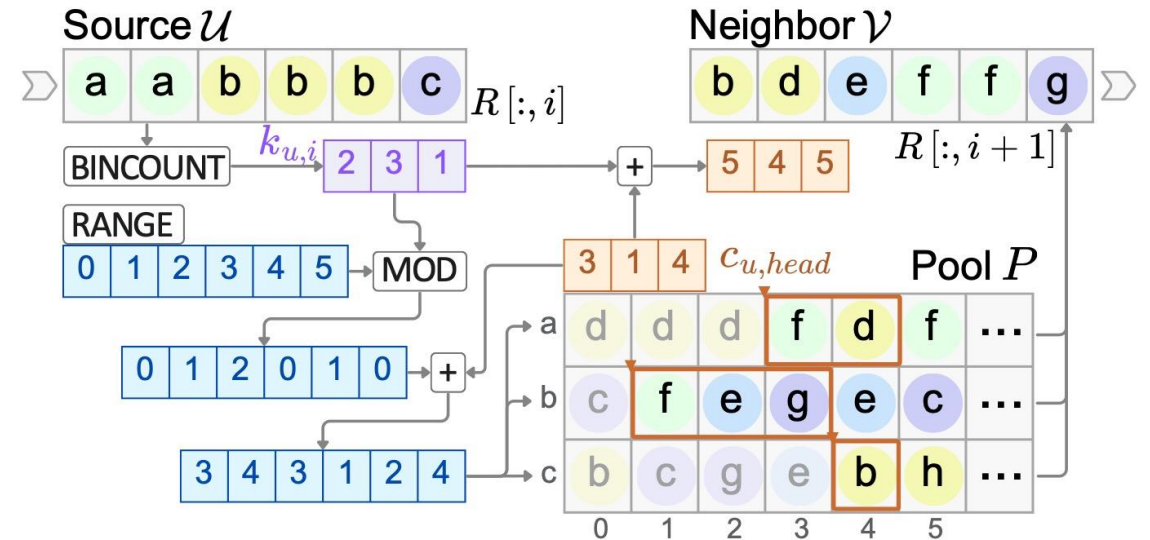
Time-based Storage

- Fast sampling across buckets
- Constrain time by sliding window



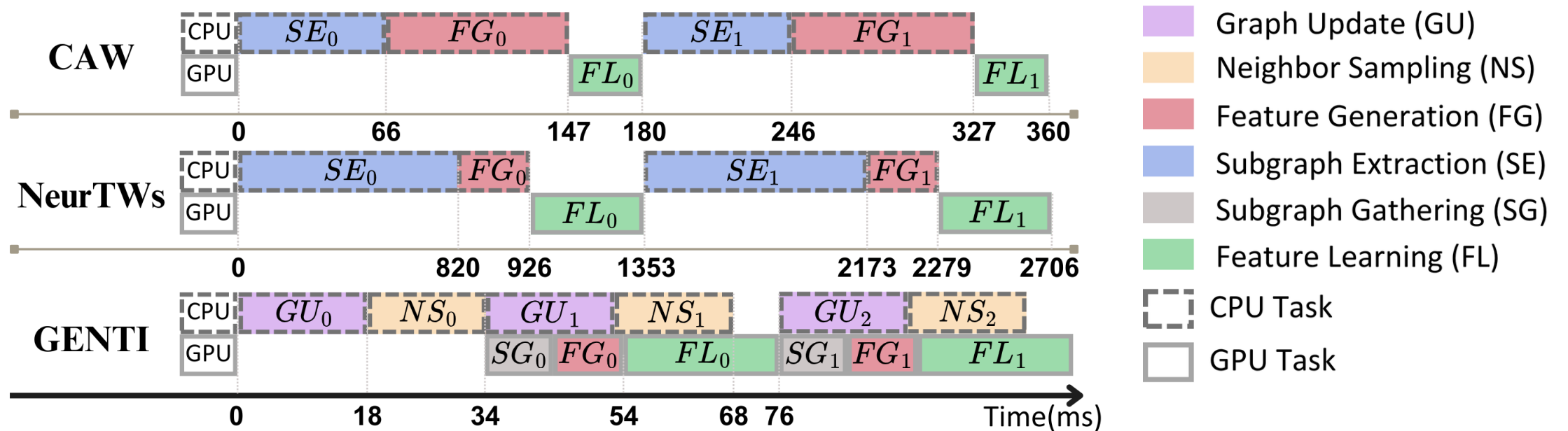
Batch Join GPU Operator

- Dense operations on GPU
- Complexity: reduce from L to \sqrt{L}



GENTI: Experimental Results

- Pipeline: high device utilization, balanced workload
- Efficiency: 3-26× speedup, fast extraction & gathering



Takeaway

HubGT: Hub Labeling for Graph Transformer

- Design: Decouple positional encoding
- Algorithm: Hierarchical hub labeling
- Experiment: Minibatch & improved accuracy

GENTI: Streaming for Subgraph GNN

- Design: Decouple subgraph extraction
- Algorithm: Batch operators & data structure on GPU
- Experiment: Tight pipeline, faster subgraph extraction

Outline

Introduction



Algorithms



Data Structures



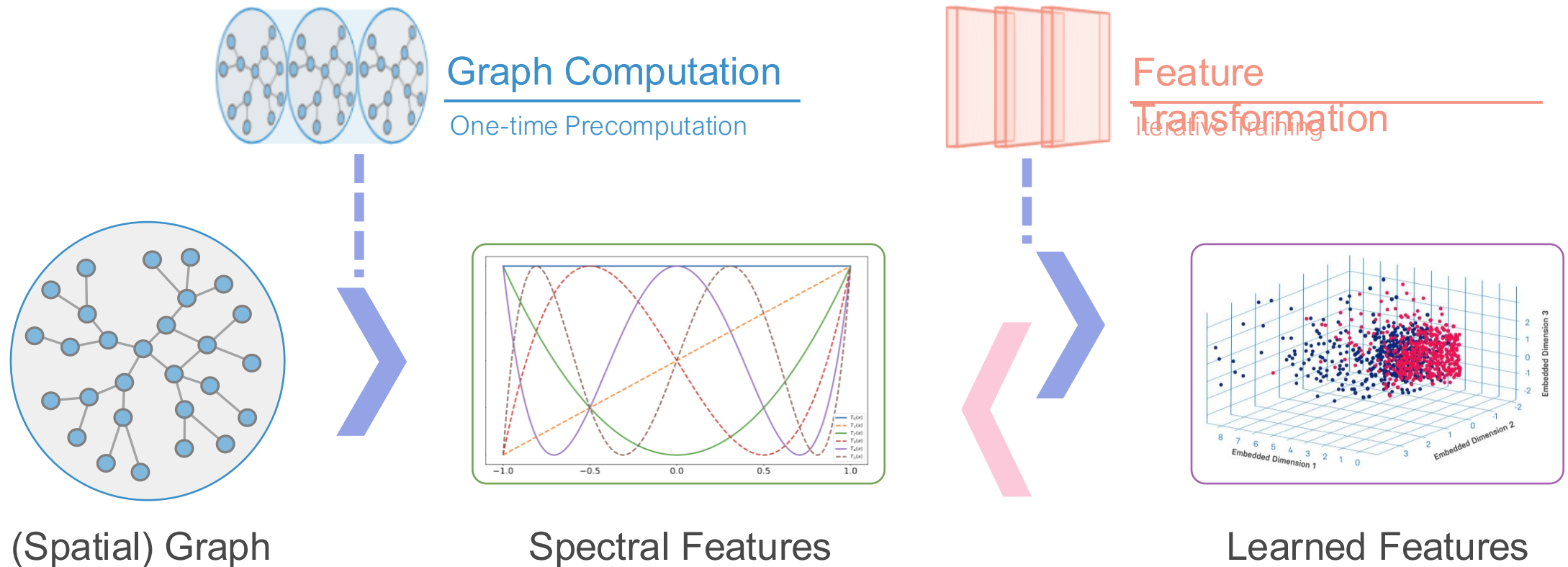
**Evaluate Scalability
Performance**

UNIFEWS: Spectral Sparsification
pyg_spectral [2020]
GN Benchmark

Conclusion and Future Directions

Background: Spectral GNN

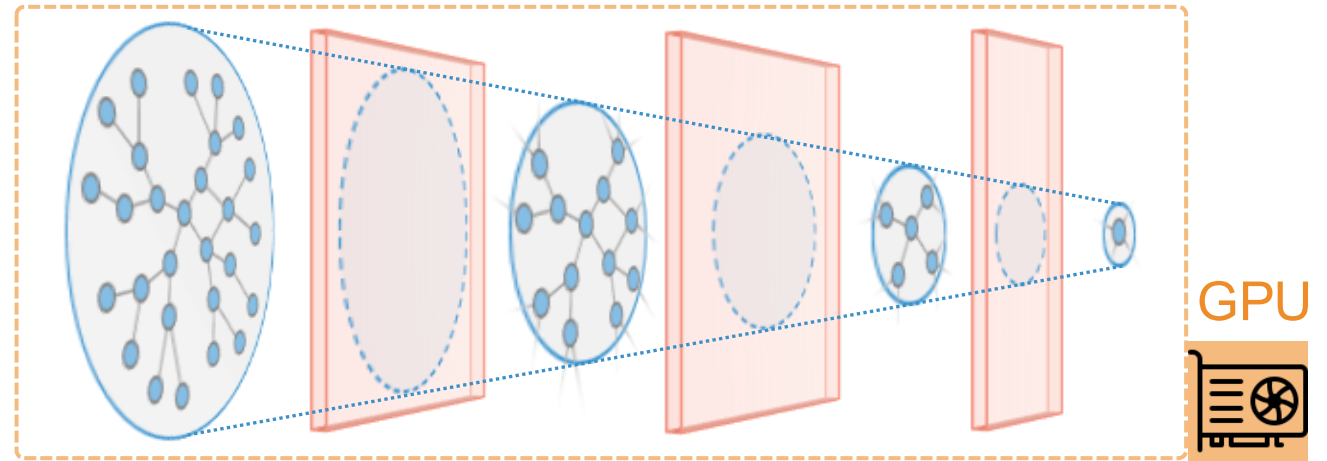
- Spectral domain: related to graph eigen-decomposition



Background: Spectral GNN

- Spatial GNN:

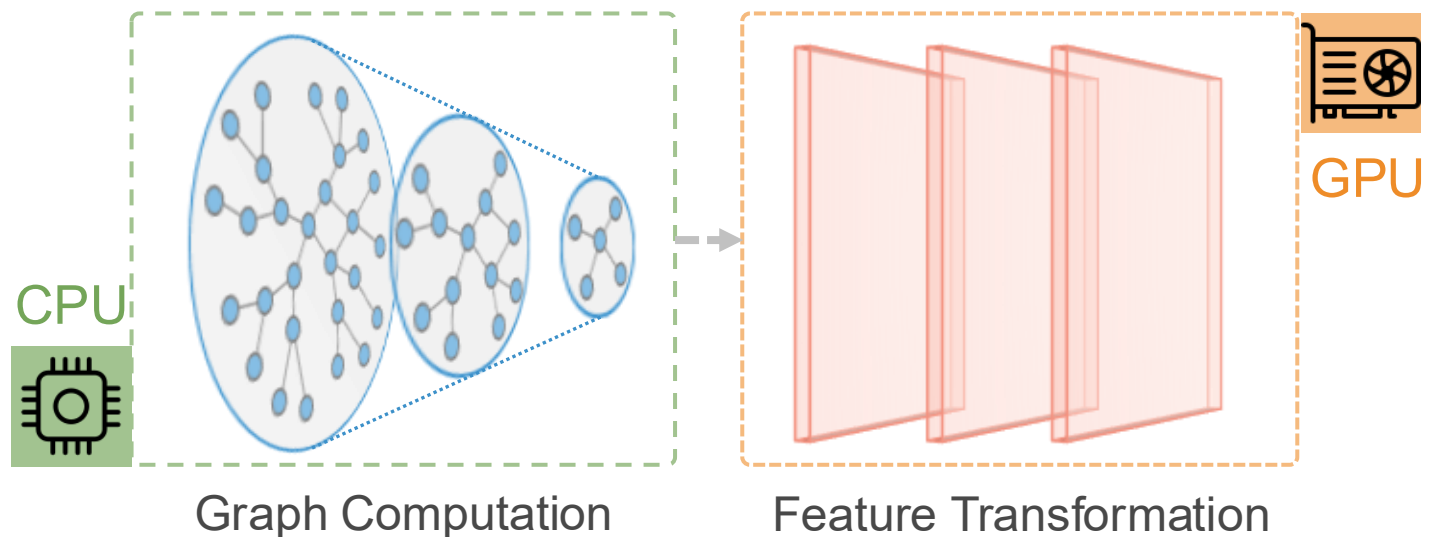
Coupled & iterative graph computation



- Spectral GNN:

Explicit & decoupled spectral features

😊 *Save memory & time*



UNIFEWS: Spectral Sparsification

- Sparsification: prune insignificant values to zero

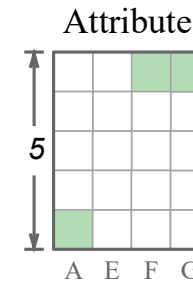
Graph Propagation

$$\mathbf{p}^{(l+1)}[u] = \sum_{v \in \mathcal{N}(u)} \mathbf{T}^{(l)}[u, v] \cdot \mathbf{p}^{(l)}[v]$$

Embedding

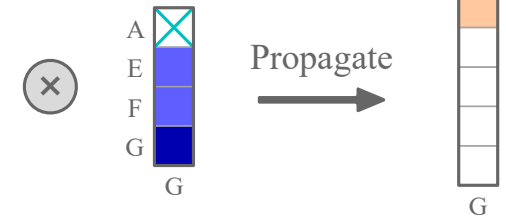
Graph Diffusion

⊗ Subject to prune



Graph Diffusion

Embedding



⊗ Pruned Entry

Feature Transformation

$$\mathbf{H}^{(l+1)}[:, i] = \sum_{j=1}^f \mathbf{W}^{(l)}[j, i] \cdot \mathbf{P}^{(l)}[:, j]$$

Representation

Weight

⊗ Subject to prune

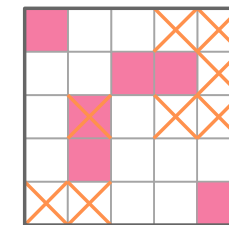
Embedding

Model Weight

Representation



⊗



Transform



⊗ Pruned Entry

UNIFEWS: Approximation Theory

- GNN as Graph Smoothing:

$$\mathbf{p}^* = \arg \min_{\mathbf{p}} \mathcal{L}, \quad \mathcal{L} = \|\mathbf{p} - \mathbf{x}\|^2 + c \cdot \mathbf{p}^\top \mathbf{L} \mathbf{p},$$

Embedding
Attribute
Graph

- Spectral Approximation:

$$\|\hat{\mathbf{p}}^* - \mathbf{p}^*\| \leq c\epsilon \|\mathbf{p}^*\|.$$

Approx Bound

- Error Guarantee:

(general)
 \Rightarrow

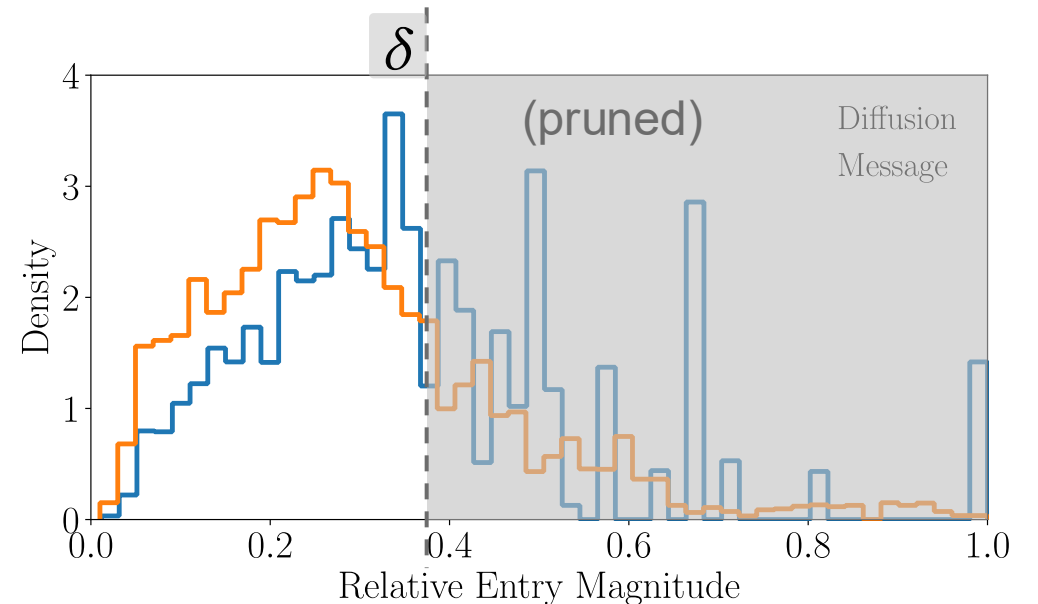
$$q_a \delta_a \leq \epsilon \|\mathbf{p}\|$$

Sparsity
Prune Threshold

(scale-free graph)
 \Rightarrow

$$\epsilon = O\left(\eta_a (1 - \eta_a)^{\frac{1}{1-\alpha}}\right)$$

Graph Distribution



UNIFEWS: Approximation Theory

Graph Sparsification:

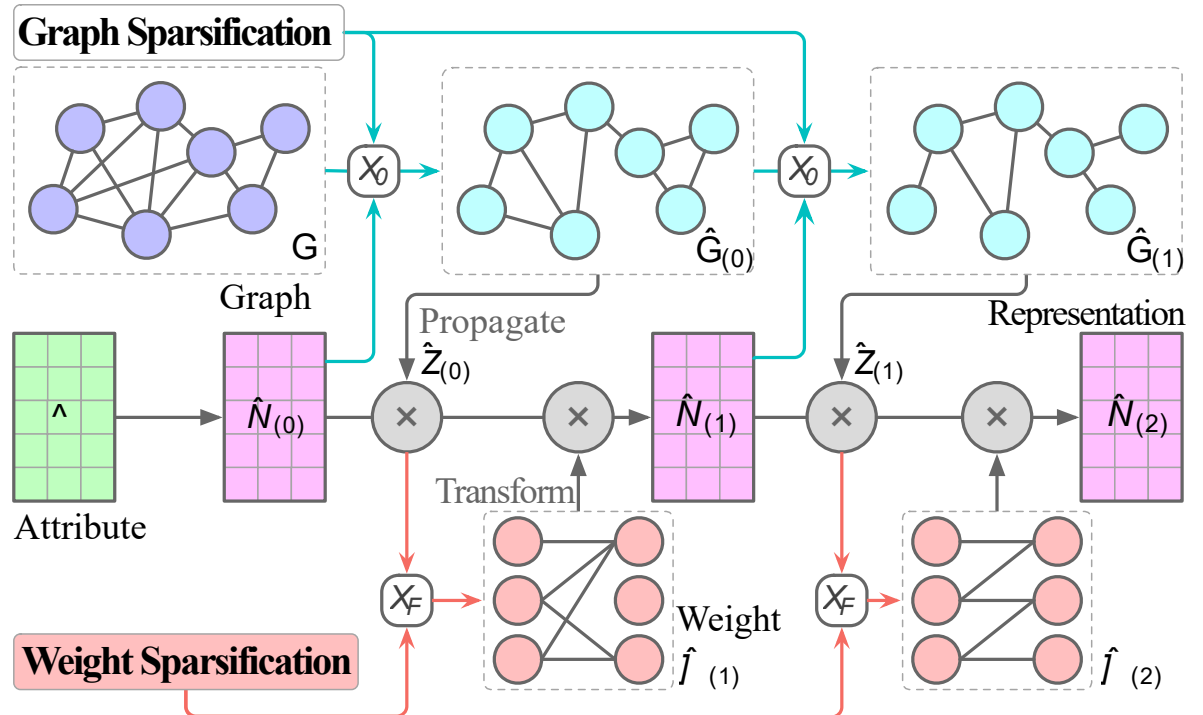
$$\| \mathbf{p}_{(l+1)} - \hat{\mathbf{p}}_{(l+1)} \| \leq O(\epsilon) \cdot \| \mathbf{p}_{(l)} \| + bc\epsilon(1 + O(\epsilon)) \| \mathbf{p}_{(l)} \|$$

Embedding / Approx Bound

+ Weight Sparsification:

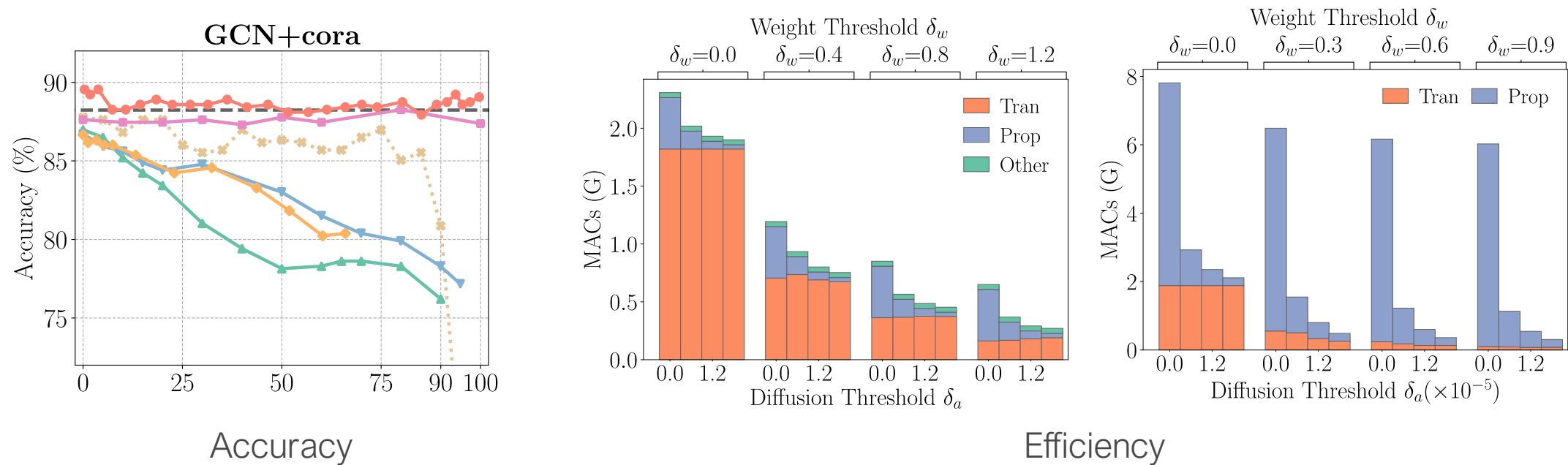
$$\| \hat{\mathbf{H}}_{(l+1)} - \mathbf{H}_{(l+1)} \|_F \leq \ell_\sigma bc\epsilon \| \mathbf{H}_{(l)} \|_F \| \mathbf{W} \|_F + \ell_\sigma q_w \delta_w$$

Representation / Network Distribution



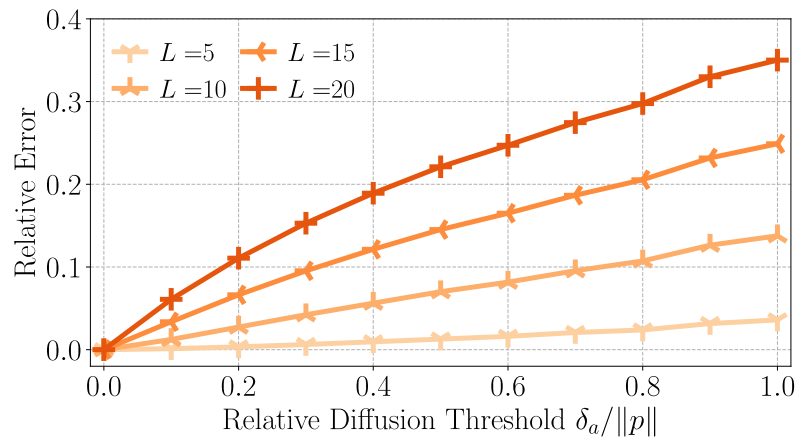
UNIFEWS: Experimental Results

- Accuracy: Accurate at high pruning ratio
- Efficiency: Boost graph & weight computation

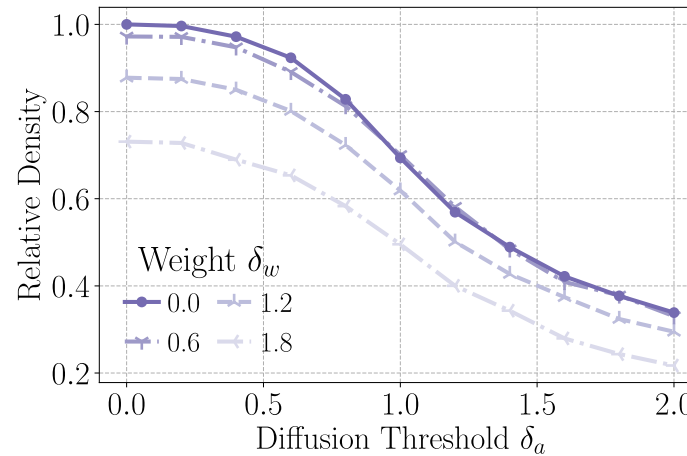


UNIFEWS: Experimental Results

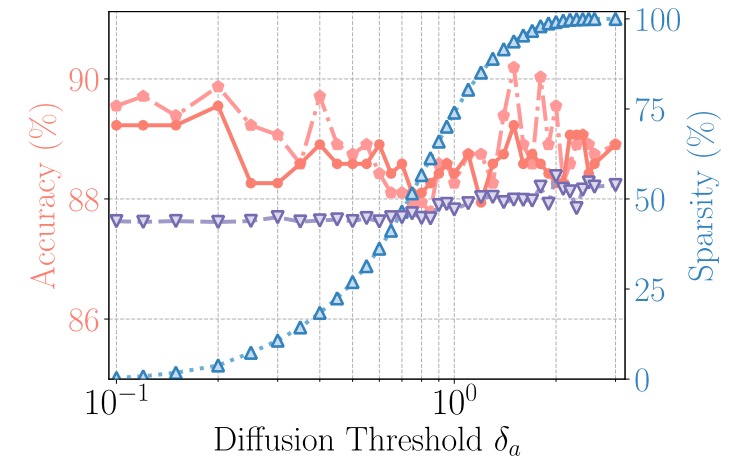
- Approximation: Follow error guarantee theory
- Implementation: Applicable to both iterative & decoupled models



Relative Error



Message Sparsity



Element Sparsity

Outline

Introduction



Algorithms



Data Structures



**Evaluate Scalability
Performance**

UNIFEWS: Spectral Sparsification

pyg_spectral: GNN Benchmark

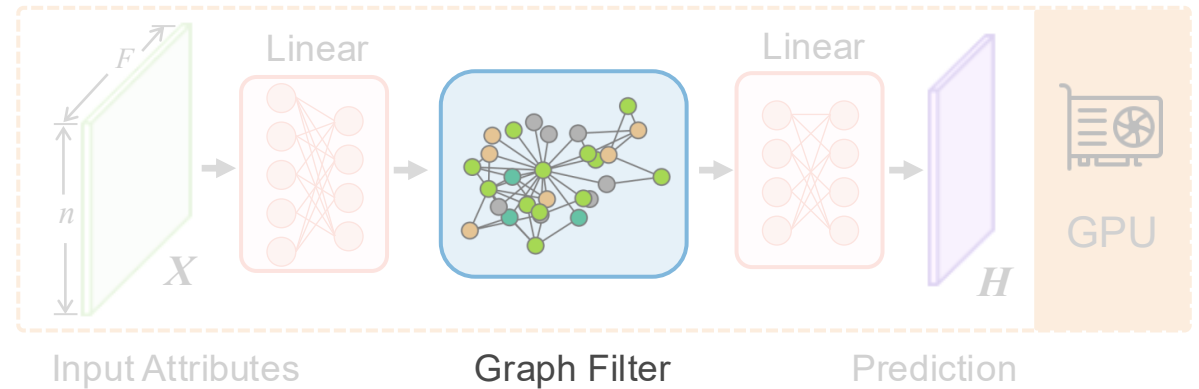
[SIGMOD 2026]

Conclusion and Future Directions

pyg_spectral: Spectral Graph Filter

- Iterative:

$$H^{(L)} = \sigma(\tilde{A}\sigma(\tilde{A}\cdots XW^{(1)}\cdots)W^{(L-1)})$$

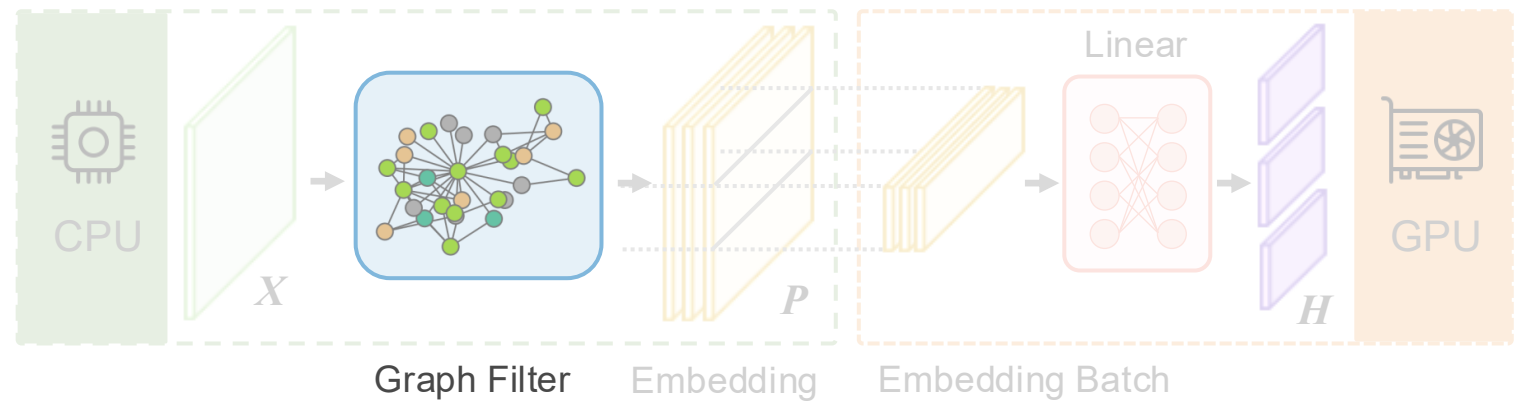


😊 *Graph filter is the same!*

- Decoupled:

$$H^{(0)} = \tilde{A}\cdots\tilde{A}X$$

$$H^{(l+1)} = \sigma(H^{(l)}W^{(l)})$$



pyg_spectral: Graph Filter Taxonomy

SCARA

$$H^{(0)} = \sum_0^{\infty} \alpha(1 - \alpha)^l \tilde{A}^l \cdot X$$

$$P_A = \phi(A^{2L} \cdot \mathcal{N})$$

LD²

$$P_{X,L2} = \sum_{l=1}^L \bar{A}^{2l} \cdot X$$

$$P_{X,H} = \sum_{l=1}^L \tilde{L}^l \cdot X$$

UNIFEWS

$$p^{(l+1)}[u] = \sum_{v \in \mathcal{N}(u)} T^{(l)}[u, v] \cdot p^{(l)}[v]$$

Spectral Filter:

$$g(\tilde{L}) = \sum_{k=0}^K \theta_k T^{(k)}(\tilde{L})$$

Filter / Graph

pyg_spectral: Graph Filter Taxonomy

GCN

$$H = (I + \tilde{A})^K \cdot X$$

APPNP

$$H = \sum_{k=0}^K \alpha(1 - \alpha)^k \tilde{A}^k \cdot X$$

GIN

$$H = (\theta I + \tilde{A})^K \cdot X$$

ChebNet

$$H = \sum_{k=0}^K \theta_k T_{\text{Cheb}}^{(k)}(\tilde{A}) \cdot X$$

BernNet

$$H = \sum_{k=0}^K \frac{\theta_k}{2^K} \binom{K}{k} (2I - \tilde{L})^{K-k} \tilde{L}^k \cdot X$$

AdaGNN

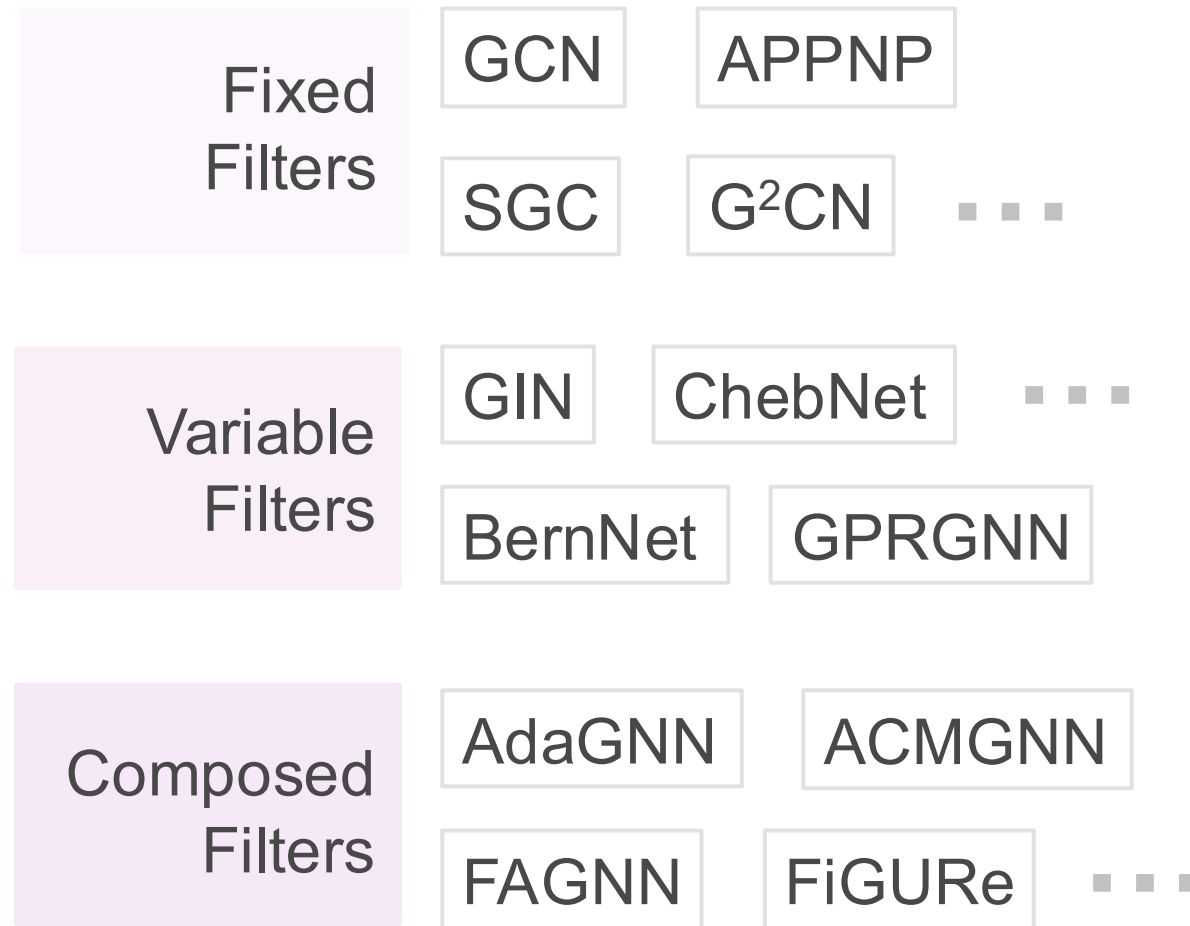
$$H = (I + \Theta_k \tilde{A})^K \cdot X$$

Spectral Filter:

$$g(\tilde{L}) = \sum_{k=0}^K \theta_k T^{(k)}(\tilde{L})$$

Filter / Graph

pyg_spectral: Graph Filter Taxonomy



Spectral Filter:

$$g(\tilde{\mathbf{L}}) = \sum_{k=0}^K \theta_k T^{(k)}(\tilde{\mathbf{L}})$$

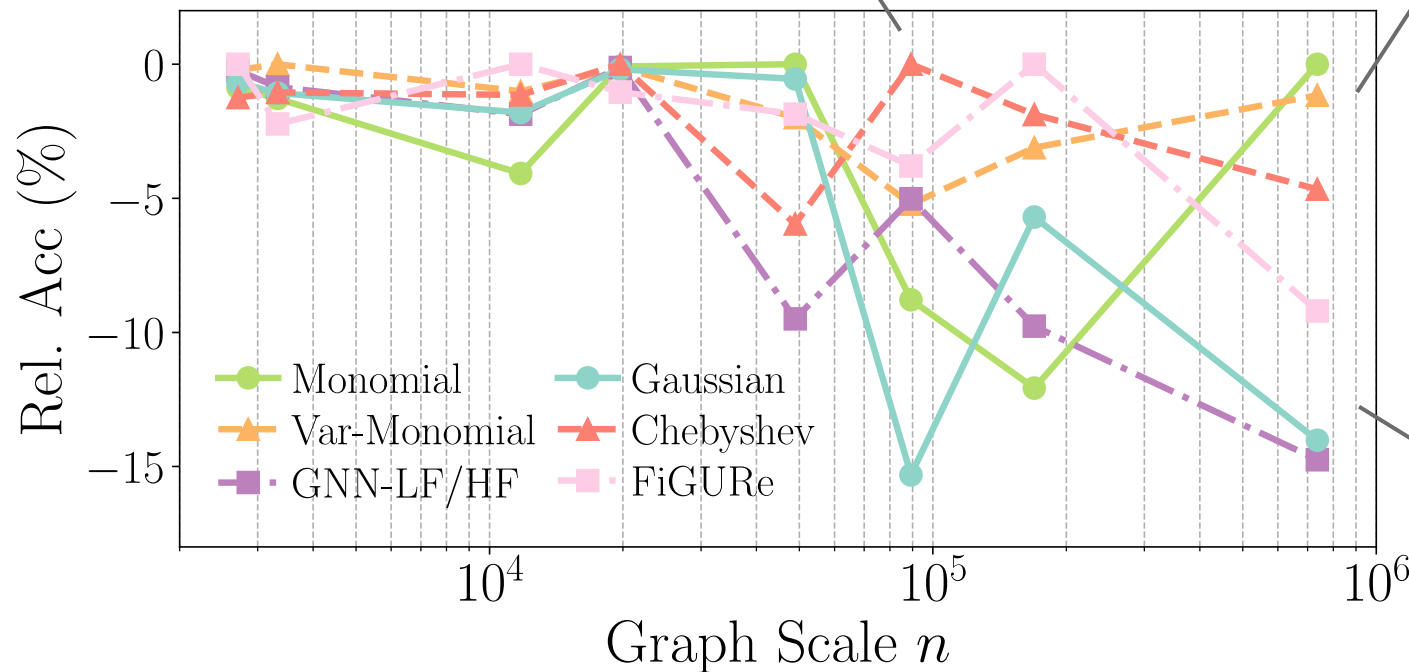
35 Covered GNNs

3 Filter Types

pyg_spectral: Key Observations

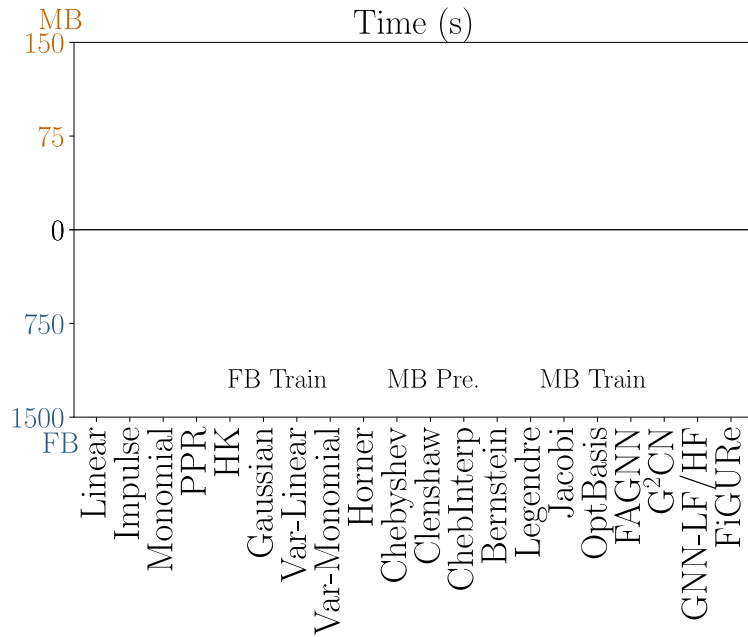
Efficacy and efficiency are not mutually exclusive.

Simple, invariable filters can also capture global view.



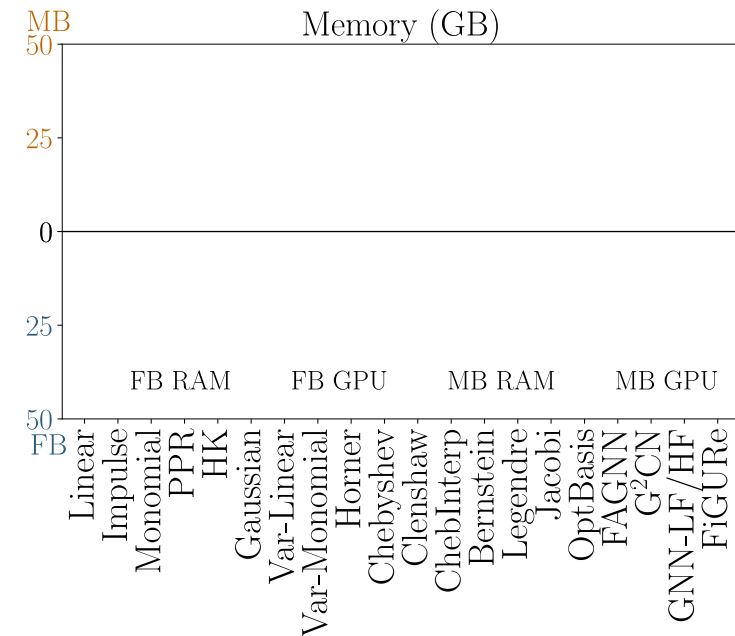
Global view is important for large-scale prediction.

pyg_spectral: Key Observations



Time Efficiency: graph operation bottleneck shifts across data scale.

Memory Overhead: Mini-batch training facilitates device utilization.



Takeaway

UNIFEWS: Spectral Sparsification Theory

- Framework: Unify graph & weight operators
- Theory: Bound progressive representation
- Experiment: Maintain accuracy & boost computation

pyg_spectral: Spectral Filter Benchmark

- Framework: Unify iterative & decoupled graph computation
- Benchmark: Evaluate time efficiency, memory footprint, & accuracy
- Insight: Harmonize effectiveness and efficiency

Outline

Introduction



Algorithms



Data Structures



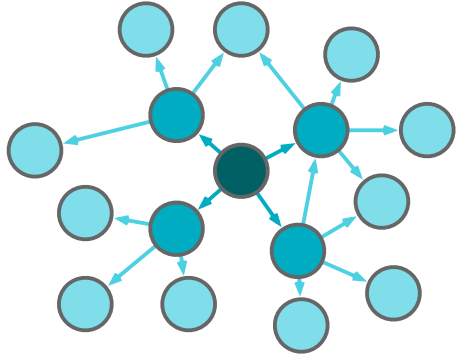
Evaluations

Conclusion and Future Directions

Scalability Challenges and Solutions

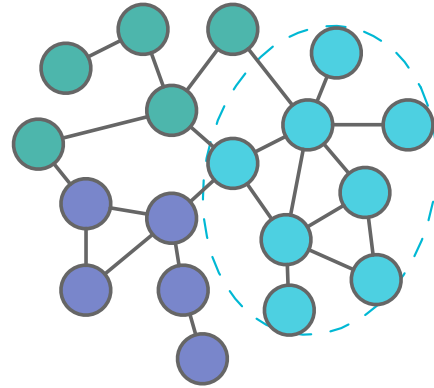
Future Research Directions

Conclusion: Scalability Challenges



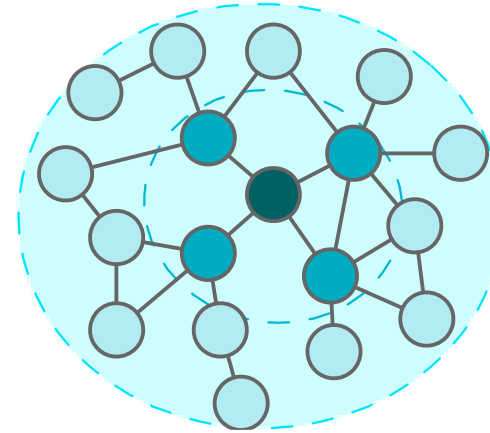
Neighbor Explosion

- Speed
- Expressive



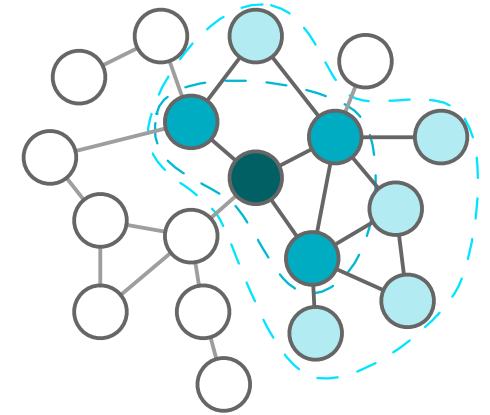
Structure Storage

- Memory
- Minibatch



Long-range

- Speed
- Heterophily
- Expressive



Subgraph

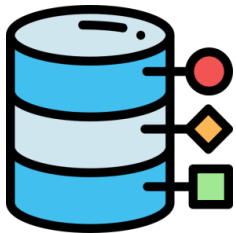
- Speed
- Fine-grained Computation

Conclusion: Scalability Solutions



Design Scalable Algorithms

- Feature-dimension Utilization [VLDB 2022 & VLDBJ 2023]
- Long-distance Embeddings [NeurIPS 2023]



Deploy Scalable Data Structures

- Hub Labeling for Graph Transformer [NeurIPS 2025]
- Streaming for Subgraph GNN [VLDB 2024]

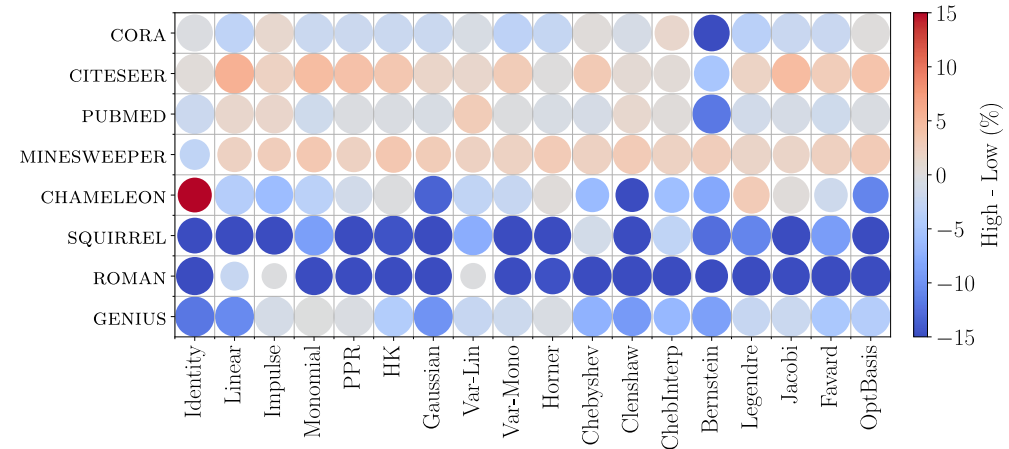
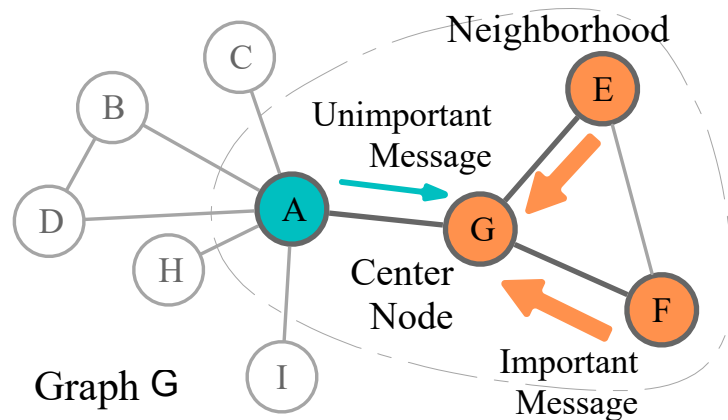


Evaluate Scalability Performance

- Spectral Sparsification Theory [ICML 2025]
- Spectral Filter Benchmark [SIGMOD 2026]

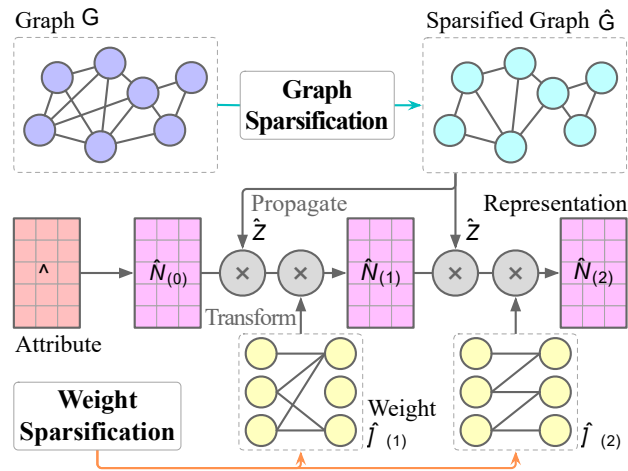
Future: Personalized Graph Analysis

- Biased performance
 - *How to ensure graph learning outcome on specific nodes?*



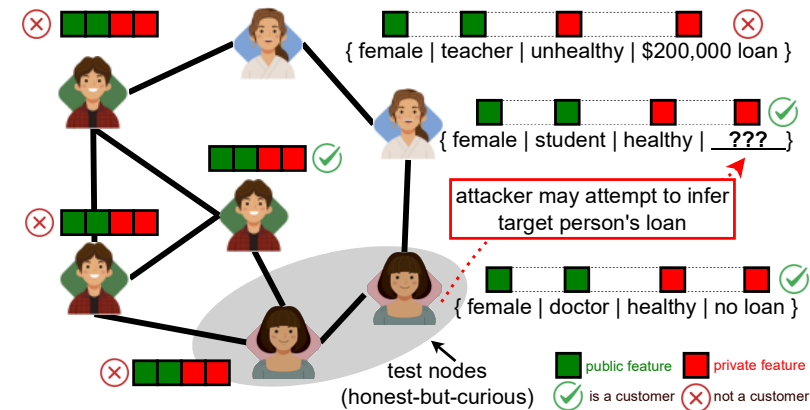
- Personalized scalability
 - *How to perform fast and scalable graph learning on specific nodes?*

Future: GNN under Various Requirements



- Computational device
 - *How to deploy GNN training/ inference on different devices?*

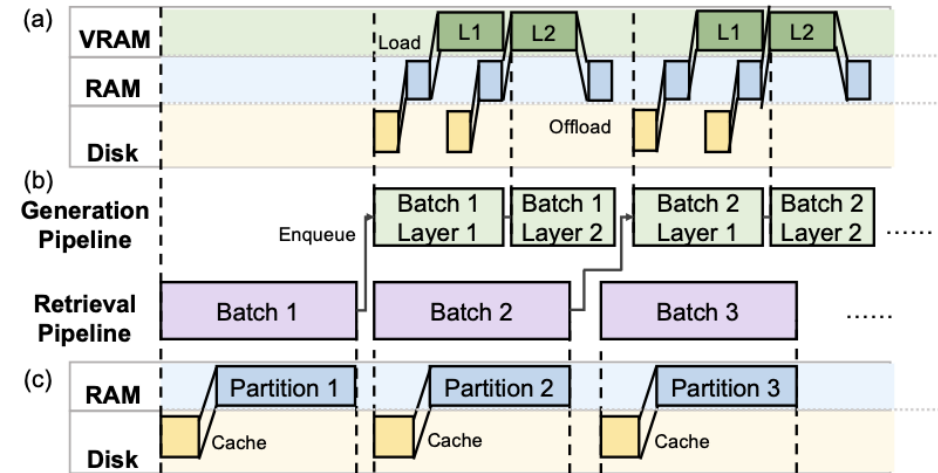
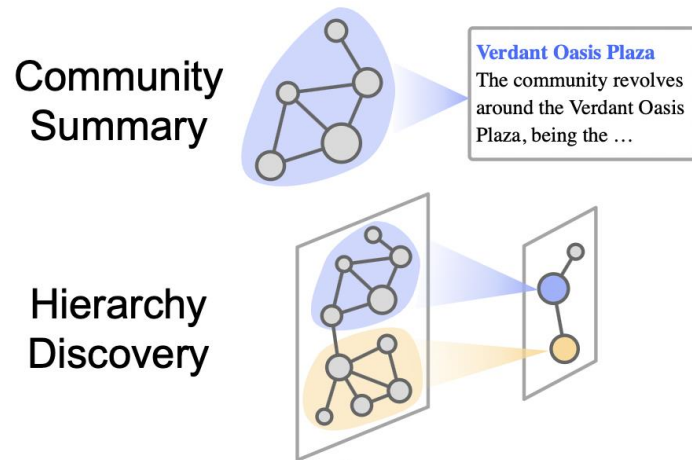
- Privacy & fairness
 - *How to effective protect local graph & attribute information?*



Future: Graph Learning + LLM

- GNN + RAG

- *How to perform efficient RAG retrieval with GNN embeddings?*



- Locality & explainability

- *How to utilize local graph structures for explainable LLM inference?*

[13] W Yu*, N Liao*, S Luo, J Liu. "RAGDoll: Efficient Offloading-based Online RAG System on a Single GPU". arXiv:2504.15302.

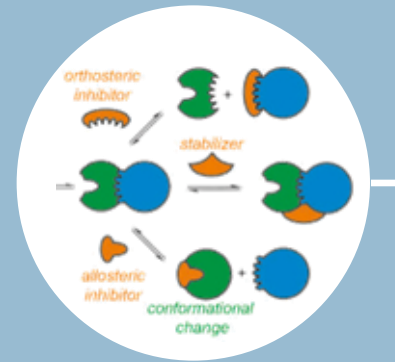
[15] I Lim*, N Liao*, Y Yang, G Yip, S Luo. "SEQ-GPT: LLM-assisted Spatial Query via Example". arXiv:2508.10486.

Vision: Towards Scalable Graph Learning

Image credits: [linkedin.com](https://www.linkedin.com), [yerliteknoloji.net](https://www.yerliteknoloji.net).



Friend Suggestion



Drug Discovery



Facial Recognition

Outline

Introduction



Algorithms



Data Structures



Evaluations

Conclusion and Future Directions

Q&A

Publications

- [1] N Liao^{*}, D Mo^{*}, S Luo, X Li, P Yin. “SCARA: Scalable Graph Neural Networks with Feature-Oriented Optimization”. **VLDB 2022**.
- [2] N Liao, D Mo, S Luo, X Li, P Yin. “Scalable Decoupling Graph Neural Network with Feature-Oriented Optimization”. **VLDBJ 2023**.
- [3] N Liao, S Luo, X Li, J Shi. “LD²: Scalable Heterophilous Graph Neural Network with Decoupled Embedding”. **NeurIPS 2023**.
- [4] N Liao^{*}, Z Yu^{*}, S Luo. “HubGT: Fast Graph Transformer with Decoupled Hierarchy Labeling”. **NeurIPS 2025**.
- [5] Z Yu^{*}, N Liao^{*}, S Luo. “GENTI: GPU-powered Walk-based Subgraph Extraction for Scalable Representation Learning on Dynamic Graphs”. **VLDB 2024**.
- [6] N Liao, Z Yu, R Zeng, S Luo. “UNIFEWS: You Need Fewer Operations for Efficient Graph Neural Network”. **ICML 2025**.
- [7] N Liao, H Liu, Z Zhu, S Luo, L Lakshmanan. “A Comprehensive Benchmark on Spectral GNNs: The Impact on Efficiency, Memory, and Effectiveness”. **SIGMOD 2026**.

Publications

- [8] H Liu, N Liao, S Luo. “Similarity-based Efficient Global Aggregation for Heterophilous Graph Neural Networks”. **ICDE 2025**.
- [9] J Yew, N Liao, D Mo, S Luo. “Example Searcher: A Spatial Query System via Example”. **ICDE Demo 2023**.
- [10] K Yow, N Liao, S Luo, R Cheng. “Machine Learning for Subgraph Extraction: Methods, Applications and Challenges”. **VLDB Tutorial 2023**.
- [11] N Liao, S Luo, X Xiao, R Cheng. “Advances in Designing Scalable Graph Neural Networks: The Perspective of Graph Data Management”. **SIGMOD Tutorial 2025**.
- [12] Y Qi, N Liao, S Luo, X Lin, J Li, K Han. “Fine-Grained Differentially Private Graph Neural Network with Node and Layer Noise Adaptation”. Under review.
- [13] W Yu*, N Liao*, S Luo, J Liu. “RAGDoll: Efficient Offloading-based Online RAG System on a Single GPU”. arXiv:2504.15302.
- [14] S Zhong, N Liao, D Mo, S Luo. “Efficient Approximate Nearest Neighbor Search with Regular Expression Filtering”. Under review.
- [15] I Lim*, N Liao*, Y Yang, G Yip, S Luo. “SEQ-GPT: LLM-assisted Spatial Query via Example”. arXiv:2508.10486.

Acknowledgements

PHD SUPERVISOR: Prof. Siqiang Luo

TAC: Prof. Cheng Long, Prof. Guosheng Lin

CO-AUTHOR/ Prof. Gao Cong (NTU), Prof. Xiaokui Xiao (NUS), Prof. Laks V.S.

COLLABORATOR: Lakshmanan (UBC), Prof. Reynold Cheng (HKU), Prof. Xiang Li (ECNU), Prof. Jieming Shi (PolyU), Dr. Pengcheng Yin (Deepmind), Dr. Cheng Chen (ByteDance)

COLLEAGUE/ Dingheng Mo, Haoyu Liu, Zihao Yu, Zhulun Zhu, Weiping Yu,

MENTEE: Shurui Zhong, Dr. Kai Wang, Dr. Yuxin Qi, Dr. Kai Siong Yow, Jun Xuan Yew, Ruixiao Zeng, Ivan Lim, Gerald Yip

FUNDING AGENCY: NTU-WeBank Joint Research Centre