Advances in Designing Scalable GNNs: The Perspective of Graph Data Management

Ningyi Liao

Siqiang Luo

Xiaokui Xiao

Reynold Cheng







2025 International Conference on Management of Data Berlin, Germany – June 22, 2025

Tutorial Contributors / Presenters



Siqiang Luo Assistant Professor, NTU

- 🖂 siqiang.luo@ntu.edu.sg
- https://siqiangluo.com
- Graph algorithms/systems, KV systems, ML data management



Ningyi Liao Ph.D Candidate, NTU

- 🖂 liao0090@e.ntu.edu.sg
- https://nyliao.github.io
- ML graph algorithms, Graph neural network



Reynold Cheng Professor, HKU

- 🖂 ckcheng@cs.hku.hk
- https://www.reynold.hku.hk
- Data science, Big graph analytics, Uncertain data management



Xiaokui Xiao

Professor, NUS

🔀 xkxiao@nus.edu.sg

https://www.comp.nus.edu.sg/~xiaoxk/

Privacy-preserving data management, Large data analysis

9:00 - 10:30 Outline

25 min 9:00 – 9:25 Introduction: Application | Challenge | Opportunity Sigiang Luo

25 min 9:25 – 9:50 **Classical Approaches:** Sampling | Partition | Decoupling Siqiang Luo & Ningyi Liao

20 min 9:50 – 10:10

20 min 10:10 – 10:30

30 min 10:30 – 11:00 Hardware-aware Sampling

Ningyi Liao

Subgraph Extraction

Ningyi Liao

Q&A Coffee Breaks 11:00 – 12:30 **Outline**

25 min 11:00 - 11:25

20 min 11:25 - 11:45

Node-wise Similarity

Sigiang Luo

Spectral Embedding Ningyi Liao

20 min 11:45 - 12:05

Graph Simplification

Ningyi Liao

25 min 12:05 - 12:30

12:30

Future Direction: Graph+LM | Specific Data | Graph System Reynold Cheng

Q&A Lunch

SIGMOD

5% SIGIVIT

GNN: Current and Prospect





Ranking in keywords of accepted papers



GNN: Learning on Graph Data

- Integrate graph computation and neural network architecture
- From semi-unstructured graph data to structured embeddings



GNN: Learning on Graph Data

Full-batch (FB) Training

- NN learns graph data as a whole
- Preserve full graph information

Mini-batch (MB) Training

- NN learns part of graph each time
- Batch division affects performance



GNN: GCN [ICLR'17]

 Message Passing: aggregate neighbor information N(u) by learnable weights W to each node as node representation h



9

$GNN: GCN \ [ICLR'17]$

• GNN Layers: Stacking multi-hop graph convolutions



T Kipf & M Welling. "Semi-supervised classification with graph convolutional networks". ICLR'17.

GNN: GAT [ICLR'18]

• Graph Attention: learnable weights for edge aggregation



GNN: Tasks and Applications



GNN: Tasks and Applications

Node-level

individuals in the network







Friend Suggestion Social Network



Knowledge Discovery Knowledge Graph

Edge-level

connections in the network

Graph-level

properties of the network





Policy Optimization Wireless Network

The Challenge: Graphs at Scale

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





The Challenge: Graphs at Scale

• Large scale graphs: diverse domains, diverse patterns

Name	Graph Size	Data Size	Domain
MAG	240M	200 GB	Heterogeneous academic graph
WikiKG	90M	160 GB	Wikidata knowledge graph
PCQM4Mv2	3M × 15	10 GB	Molecules in quantum chemistry interaction
Amazon	2M	8 GB	Product co-purchasing network
MD17	3M × 15	8 GB	Molecular dynamics trajectories
COCO-SP	120K × 500	12 GB	Segmentation of computer vision images
Snap-Patents	3M	20 GB	Patent citation network
Hyperlink Graph	4B	400 GB	Crawled webpage hyperlink

The Challenge: 📽 Neighbor Explosion

- Neighboring nodes intensively increases with network depth
- Exponential, related with edge size
- Lost of node identity





The Challenge: 📽 Neighbor Explosion

Case study: GCN computation graph



The Challenge: # Limited Memory

- GPU memory is limited and costly
- Storing and processing graph data requires large memory space





The Challenge: SMulti-scale

- Not only local nodes are important in prediction
- Need to capture global graph information efficiently
- <u>Homophily</u>: the inductive case, similar nodes cluster together
- <u>Heterophily</u>: the opposite case, connected nodes tend to be of dissimilar classes



The Challenge: SMulti-scale

Example: Heterophily in fraud detection

- Graph: financial transaction network with benign or fraud users
- Task: identify fraud users
- Fraudsters use benign accounts as springboards sto reach other fraudsters







The Opportunity: Graph Data Management



What can the DB community do for GNN challenges?



Preview: 📽 Neighbor Explosion

Solution from Node-wise Similarity: GBP [NeurIPS'20]

- Replace GNN computation by Personalized PageRank
- Reduce neighborhood by bidirectional search
- Decoupling CPU and GPU computation



Preview: # Limited Memory

Solution from Subgraph Extraction: GENTI [VLDB'24]

- Only load subgraph (part of the graph) for GNN computation
- Structured subgraph sampled by *k*-many random walks
- GPU-oriented data structure and batch computation



Preview: SMulti-scale

Solution from Spectral Embedding: LD² [NeurIPS'23]

- Embed full-graph information in spectral domain
- Bridge spatial and spectral longrange graph operations
- Fast & scalable spectral embedding computation



Preview: Scalable Graph+LM

Future Direction in Graph-based RAG

- Scalable graph embedding, indexing, and retrieval
- Integration between graph data & LLM



Advances in Designing Scalable GNNs

Classical Approaches for Scalable GNNs

Outline

Introduction

25 min 9:25 – 9:50 Classical Approaches for Scalable GNNs Siqiang Luo & Ningyi Liao

Basic GNN: GCN (T Kipf & M Welling | ICLR'17)
Graph Sampling: GraphSAGE (W Hamilton et al | NeurIPS'17)
Graph Partition: AliGraph (R Zhu et al | VLDB'19)
Decoupled Propagation: SGC (F Wu et al | ICLR'19)
Application: Social Network Friend Recommendation
Hardware-aware Sampling

Subgraph Extraction

Basic GNN: GCN Scalability Analysis



Basic GNN: GCN Scalability Analysis

Graph-wise view

Node-wise view

Computation



Basic GNN: GCN Complexity

Node-wise view

Computation



Graph Sampling

- <u>Node-wise</u> sampling: limit neighborhood size of each node
- <u>Layer-wise</u> sampling: limit total neighborhood size of each layer
- <u>Subgraph-wise</u> sampling: limit neighborhood within a subgraph



X Liu et al, "Sampling Methods for Efficient Training of Graph Convolutional Networks: A Survey", *J. Autom. Sinica 2022*.



Graph Sampling: GraphSAGE (NeurIPS'17)

- Sampling in graph aggregation: limit neighborhood size of each node to *S*
- Structured neighborhood: sampled representations are fixed-length
- Feasible for various aggregation schemes





Classical Approaches

S = 2 Computation Graph

W Hamilton et al, "Inductive Representation Learning in Large Attributed Graphs", NeurIPS'17.

Graph Partition

- Objective: minimize data loss in subgraphs during training
- Commonly used in multi-host/ multi-GPU training
- MB: only load subgraph onto GPU





S Bajaj et al, "Graph Neural Network Training Systems: A Performance Comparison of Full-Graph and Mini-Batch", *VLDB'25*.

Graph Partition: AliGraph (VLDB'19)

- Available partitions:
 - METIS
 - Vertex/Edge cut
 - 2D partition
 - Streaming
- Local Sampling within subgraphs on each host



Classical Approaches

Decoupled Propagation

- Graph propagation is the computational bottleneck
- Graph computation \rightarrow CPU
- Network learning \rightarrow GPU
- MB: random sampling/ hashing



S Bajaj et al, "Graph Neural Network Training Systems: A Performance Comparison of Full-Graph and Mini-Batch", *VLDB'25.* Full-batch



Mini-batch

CPU+GPU Workload



N Liao et al, "A Comprehensive Benchmark on Spectral GNNs: The Impact on Efficiency, Memory, and Effectiveness", *SIGMOD'26*.

Decoupled Propagation: SGC (ICLR'19)

- Combine the iterative propagations in GCN
- Reduce GPU memory overhead: $O(m) \rightarrow O(s)$
- No change in time overhead: remain O(Lm)

GCN (iterative)



SGC (decoupled)



Classical Approaches: Takeaways

- Graph computation is the bottleneck in GNN scalability
- Challenges:





• Approaches:

Graph Sampling

• Reduce neighborhood aggregation

Graph Partition

• Reduce subgraph size

Graph Partition

• Batch training on subgraphs

Decoupled Propagation

• Transfer computation to CPU
Classical Approaches: Applications

Social Network Friend Recommendation

- Candidates are partitioned: based on different retrieval criteria
- GraphSAGE sampling: fast neighbor
 lookup
- MB: sampling within partitions



Classical Approaches: Evaluation

Memory Footprint

- Sampling: suffer from neighbor explosion for using more redundant memory
- **Decoupling**: better memory scalability for less additional memory overhead
- Possible OOM under limited memory



N Liao et al, "Scalable Decoupling Graph Neural Networks with Feature-Oriented Optimization", VLDBJ'23.



Classical Approaches: Evaluation

Time Efficiency

- Sampling: higher time overhead caused by device idleness and transmission
- **Decoupling**: speed mostly determined by graph computation algorithms





N Liao et al, "Scalable Decoupling Graph Neural Networks with Feature-Oriented Optimization", VLDBJ'23.

Advances in Designing Scalable GNNs

Hardware-aware Sampling

Outline

Introduction

Classical Approaches

20 min 9:50 – 10:10

Hardware-aware Sampling

) Ningyi Liao

TransmissionADGNN [Z Song et al | SIGMOD'23]GIDS [J Park et al | VLDB'24]ComputationNeutronOrch [X Ai et al | VLDB'24]DAHA [Z Li et al | VLDB'24]

Subgraph Extraction

Hardware-aware Sampling

Transmission: ADGNN [SIGMOD'23]



Global Sampling

- Full graph view
- Precision guarantee
- Transmission overhead



Local Sampling

- Distributed storage & computation
- Information loss
- Depends on partition

Transmission: ADGNN [SIGMOD'23]

Balance Local & Global Sampling

- Maintain a set of effective neighbors
 locally
- Only communicate for top remote nodes that reduce precision





Hardware-aware Sampling



Hardware-aware Sampling

Transmission: GIDS [MLDB'24]



Graph Structure: Parallel Sampling

- Sampling can be performed independently from NN training
- Parallel sampling for batches on CPU graph storage



J Yang et al, "GNNLab: a factored system for sample-based GNN training over GPUs", *EuroSys'22.*

J Park et al, "Accelerating Sampling and Aggregation Operations in GNN Frameworks with GPU Initiated Direct Storage Accesses", VLDB'24.

Transmission: GIDS [VLDB'24]



Hardware-aware Sampling

Node Attribute: Hierarchical Buffer



J Park et al, "Accelerating Sampling and Aggregation Operations in GNN Frameworks with GPU Initiated Direct Storage Accesses", VLDB'24.

Computation: NeutronOrch [VLDB'24]

Layer-based Computation Orchestrating

- CPU: 1st layer, from attribute to embedding
- GPU: higher layers, only process embeddings





SGT

Feature transfer

Hardware-aware Sampling



X Ai et al, "NeutronOrch: Rethinking Sample-Based GNN Training under CPU-GPU Heterogeneous Environments", *VLDB'24*.

Hardware-aware Sampling

Computation: NeutronOrch [VLDB'24]

Cross-batch Data Pipelining

- Hotness-aware sampling across batches on CPU
- Constrain staleness within super-batch



Computation: DAHA [VLDB'24]

尊

Hardware-aware Sampling

Profiling & scheduling between CPU and GPU computation



Hardware-aware Sampling

14/

1.1



and Hardware Aware Execution Planning", VLDB'24.

GPU

Z=AT





Z Li et al, "DAHA: Accelerating GNN Training with Data and Hardware Aware Execution Planning", *VLDB'24*.



Hardware-aware Sampling: Takeaways

Transmission

Computation

ADGNN [SIGMOD'23]

- Between computational devices
- 。Balance Local & Global Sampling

GIDS [VLDB'24]

Between storage devices
Hierarchical node attribute buffer

NeutronOrch [VLDB'24]

- Among GNN operations
- Layer- & batch-wise pipelining

DAHA [VLDB'24]

- Among data partitions
- 。 Balance CPU & GPU computation

Advances in Designing Scalable GNNs

Subgraph Extraction

Outline

Introduction

Classical Approaches

Hardware-aware Sampling

20 min 10:10 – 10:30

Subgraph Extraction

Ningyi Liao

Partition G3 [X Wan et al | SIGMOD'23]

Generation SUREL [H Yin et al | VLDB'23] & GENTI [Z Yu et al | VLDB'24]

Slicing TIGER (K Wang et al | VLDB'24)

Application: Hardware Obfuscation Analysis

Subgraph Extraction: Challenges

S Effectiveness

 Topology-based partition may be not optimal



Efficiency

 Data transferring overhead is significant



Partition: G3 [SIGMOD'23]



- Hybrid Parallelism
 - Data-parallel: graph partitions on workers
 - Model-parallel: full-graph training with P2P transmission
 - Pipeline-parallel: overlapped computation and transmission



X Wan et al, "Scalable and Efficient Full-Graph GNN Training for Large Graphs", *SIGMOD'23*.

Partition: G3 [SIGMOD'23]

- 2-stage graph partition
 - Stage 1: unsupervised, minimize global communication by reducing cut edges
 - Stage 2: adaptive updated, balance empirical communication



Subgraph Extraction

Generation: SUREL [VLDB'23]

- Walk-based Subgraph Extraction
 - Sample multiple walks for each query node W_u
 - Join walks to generate a subgraph



H Yin et al, "Algorithm and system so-design for efficient subgraph-based graph representation learning", VLDB'22. H Yin et al, "SUREL+: Moving from Walks to Sets for Scalable Subgraph-Based Graph Representation Learning", VLDB'23.

Subgraph Extraction

 \mathcal{A}

(0)

(1)

(2)

(3)

X

 $\mathcal{X}_{u,u}$ [2,0,0]

 $\mathcal{X}_{u,b} [0, 2, 0] \\ \mathcal{X}_{u,a} [0, 0, 1]$

 $\mathcal{X}_{u,v}$ [0, 0, 1]

 $\mathcal{X}_{v,u}$ [0,0,1]

 $\mathcal{X}_{v,b}$ [0,2,0]



 \mathcal{T}

 \mathcal{H}

 $\mathcal{H}_u: \mathcal{V}_u \to \mathcal{X}_u(\text{RPE-ID}_u)$

Generation: SUREL [MLDB'23]



Subgraph Extraction

Walk-based Storage

- Fast neighbor access with pointer
- Reduce feature overlap



Sparse Join Operator

- Remove duplicate nodes
- Reduce sparse storage



H Yin et al, "Algorithm and system co-design for efficient subgraph-based graph representation learning", *VLDB'22*. H Yin et al, "SUREL+: Moving from Walks to Sets for Scalable Subgraph-Based Graph Representation Learning", *VLDB'23*.

Generation: GENTI [VLDB'24]



• Dynamic Subgraph Extraction: streaming of subgraph sampling & join



Generation: GENTI [VLDB'24]



Time-based Storage

- Fast sampling across buckets
- Constrain time by sliding window

		κλ			
(a)	间 Stale	Active		Slic	ling Window
Neighbor	bec	f h b	ottom	g	d a top
Lookup Table	0 1 2	0 1	\overline{xt} in	ext 0	0 1
(b)					
Bucket	${\mathcal I}_{u,0}$	${\mathcal I}_{u,1}$	•••	${\mathcal I}_{u,r-1}$	${\mathcal I}_{u,r}$
Timestamp	[0, <i>α</i>)	$[\alpha, 2\alpha)$		[(r-1)lpha,rlpha)	[rlpha,(r+1)lpha)
Sample Weight	$\left[0,e^1 ight)$	$\left[e^{1},e^{2} ight)$		$\left[e^{r-1},e^r ight)$	$\left[e^{r},e^{r+1} ight)$
Random Number	Ø	$\{z_4\}$	•••	$\{z_3\}$	$\{z_1,z_2\}$
Fetch Amount	0	1		1	2
Sampled Neighbor	Ø	f		g	a d

Batch Join GPU Operator

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



Z Yu et al, "GENTI: GPU-Powered Walk-Based Subgraph Extraction for Scalable Representation Learning on Dynamic Graphs", *VLDB*'24.

Slicing: TIGER [VLDB'24]

- Subgraph extraction on heterogeneous graph: based on edge triples
- Bottleneck: element lookup on disk
- <u>Slicing</u>: store elements of same triples together ← NP-Hard



Slice-based Extraction

- Sequential disk access
- Slicing overhead

Slicing: TIGER [VLDB'24]

- Slicing: match-then-generate to reduce overlap storage
- Caching: recently accessed elements in memory





Subgraph Extraction: Takeaways

Partition	Generation	Slicing	
G3 [SIGMOD'23]	SUREL [VLDB'23] GENTI [VLDB'24]	TIGER [VLDB'24]	
 3 Parallel enhancements Communication overhead persist under full-graph training 	 Generated subgraphs may be more expressive Generation and gathering overhead 	 On-disk subgraph extraction 2-stage slicing and caching for batter local access 	

Subgraph Extraction: Evaluation

Model Operations

- Hash: constantly fastest
 & insignificant
- Metis family: fast but less scalable
- Streaming family: slow & less scalable







- Mini-batch training is more scalable than full-graph training
- Mini-batch + local sampling is effective in reducing computation





S Bajaj et al, "Graph Neural Network Training Systems: A Performance Comparison of Full-Graph and Mini-Batch", VLDB'25.

Subgraph Extraction: Applications

Hardware Obfuscation Analysis

- Graph: integrated circuits logic gates
- Task: identify obfuscated gates & interconnections



L Mankali et al, "Titan: Security Analysis of Large-Scale Hardware Obfuscation Using Graph Neural Networks", *TIFS'22*.

Outline

Q&A 30 min 10:30 - 11:00 25 min 11:00 - 11:25 20 min 11:25 - 11:45 20 min 11:45 - 12:05 25 min 12:05 - 12:30

Coffee Breaks **Node-wise Similarity**

Sigiang Luo

Spectral Embedding

Ningyi Liao

Graph Simplification

Ningyi Liao

Future Direction: Graph+LM | Specific Data | Graph System Reynold Cheng

Advances in Designing Scalable GNNs

Node-wise Similarity

Outline

25 min 11:00 – 11:25

Node-wise Similarity

Siqiang Luo

Personalized PageRankGBP [Chen et al | NeurIPS'20] & SCARA [Liao et al | VLDB'22]Pair-wise SimilaritySIGMA [H Liu et al | ICDE'25]Hub LabelingCFGNN [Preprint] | HubGT [Preprint]Application:Recommendation Popularity

Spectral Embedding

Graph Simplification

Future Direction

Node Similarity: Schema



Node Similarity

Personalized PageRank (PPR)



- A *a*-decay random walk:
 - Stop with *a* probability at each step
 - If not stop, randomly jumps to one of its out-neighbors

Node Similarity

Personalized PageRank (PPR)





- Random walk interpretation of PPR(s, t):
 - probability that a random walk from s stops at (
Node Similarity

Personalized PageRank (PPR)

Bidirectional push for faster node-wise computation



P Lofgren et al, "Personalized PageRank Estimation and Search: A Bidirectional Approach", *WSDM'16*.

PPR: GBP [NeurIPS'20]

- Bidirectional propagation
 - Forward: Random Walk from ego nodes
 - Backward: Reverse Push from feature vectors
- Decoupled propagation and network
- Reduced propagation complexity

 $O(L\sqrt{Lm\log(Ln)}/\varepsilon)$



- → Forward: Random Walk
- → Backward: Reverse Push

Node Similarity



PPR: SCARA [VLDB'22]

- <u>Feature-oriented</u>: start from each node attribute vector
- Pure forward scheme
 - Insignificant push: control by threshold
 - Significant push: compensate by random walk
- Reduced propagation complexity

 $O\left(\sqrt{m\log(n)}/\lambda\right)$

Node Similarity



N Liao et al, "SCARA: Scalable Graph Neural Networks with Feature-Oriented Optimization", *VLDB'22*. N Liao et al, "Scalable Decoupling Graph Neural Networks with Feature-Oriented Optimization", *VLDBJ'23*.

PPR: SCARA [VLDB'22]

- Conventional PPR: each node possesses a scalar value
- PPR in GNN: each node possesses a vector (feature vector)
- There are overlaps among feature dimensions



We can reuse PPR calculations among feature dimensions!

п $\hat{\boldsymbol{\pi}}(\boldsymbol{b}_i, \boldsymbol{\gamma} \boldsymbol{\beta}_s)$ Base **Base Calculation** (High Precision) (2) Non-base $\pi^*(x_f)$ $\sum \theta_i \cdot \hat{\boldsymbol{\pi}}(\boldsymbol{b}_i, \boldsymbol{\gamma} \boldsymbol{\beta}_s)$ **Base** Combination

(1)

 $\hat{\boldsymbol{\pi}}(\boldsymbol{x}', (1-\gamma \sum \theta_i) \beta_s)$ **Residue** Calculation (Low Precision)

N Liao et al, "SCARA: Scalable Graph Neural Networks with Feature-Oriented Optimization", VLDB'22. N Liao et al, "Scalable Decoupling Graph Neural Networks with Feature-Oriented Optimization", VLDBJ'23.

Node Similarity

Node Similarity

Pair-wise Similarity



- <u>SimRank</u>: a measure of node similarity based on neighborhood topology
- Global Similarity: distant node pair sharing the similar topology



Pair-wise Similarity: SIGMA [ICDE'25]

• Global similarity: address heterophily by graph metric



Node Similarity

H Liu et al, "SIGMA: An Efficient Heterophilous Graph Neural Network with Fast Global Aggregation", *ICDE'25*.

H Liu et al, "SIGMA: An Efficient Heterophilous Graph Neural Network with Fast Global Aggregation", *ICDE'25*.

Pair-wise Similarity: SIGMA [ICDE'25]

• Combining attribute and topological representations:

$$\mathbf{H} = \mathrm{MLP}_{H} \left(\delta \cdot \mathbf{H}_{\mathbf{X}} + (1 - \delta) \cdot \mathbf{H}_{\mathbf{A}} \right).$$

• Combining local and global representations:

$$\begin{array}{l} \text{AGG}: \ \widehat{\mathbf{Z}}_u = \sum_{v \in V} \mathbf{S}(u, v) \cdot \mathbf{H}_v, \\ \text{UPD}: \ \mathbf{Z}_u = (1 - \alpha) \cdot \widehat{\mathbf{Z}}_u + \alpha \cdot \mathbf{H}_u, \end{array}$$





Node Similarity

Hub Labeling

- <u>Hub Labeling</u>: adding new links to selected hub nodes, so that distant node pairs can be easily reached through hubs
- <u>Shortest Path</u>: reachability and distance can be computed through labels



Node Similarity

Hub Labeling: CFGNN [Preprint]

- Hub labeling: identify hub nodes as the center of node clusters
- Distinctive message-passing:
 - Fringe nodes -> Hub nodes: collect
 - Hub nodes -> Fringe nodes: distribute



Node Similarity

Hub Labeling: HubGT [Preprint]

- Hub Labeling: compute node-pair Shortest Path Distance (SPD)
- <u>Positional Encoding</u>: emphasize important pairs in Graph Transformer

$$ilde{oldsymbol{H}}_i = ext{softmax} \left(rac{oldsymbol{Q}_i oldsymbol{K}_i^ op}{\sqrt{F_K}} + oldsymbol{P}
ight) oldsymbol{V}_i,$$

• Complexity: $O(n^2)$ (full graph) $O(s^2)$ (within subgraph)





Hub Labeling: HubGT [Preprint]

 \Diamond

Node Similarity

• Subgraph SPD: fast querying any node pairs within subgraph of a node



N Liao et al, "HubGT: Fast Graph Transformer with Decoupled Hierarchy Labeling", *arXiv:2412.04738, 2024*.

Node Similarity: Takeaways

PPR	SimRank	Hub Labeling
GBP [NeurIPS'20] SCARA [VLDB'22]	SIGMA [ICDE'25]	CFGNN [Preprint] HubGT [Preprint]
 Similar to graph convolution with faster computation Scalar → feature vector 	 Global similarity between distant nodes Fast push computation 	 Rewiring with core hub nodes Applicable to MPNN & GT

Node Similarity: Applications

Debiasing Recommendation

- Graph: user-item interaction
- Task: unbiased recommendation
- Global popularity: biased to top popular items
- Personal popularity: based on user similarity



Advances in Designing Scalable GNNs

Spectral Embedding

Outline

Node-wise Similarity

20 min 11:25 – 11:45

Spectral Embedding

Ningyi Liao

Channel CombinationLD2 [Liao et al | NeurIPS'23]TER [Wang et al | KDD'23] & S3GCL [Wan et al | ICML'24]Adaptive SelectionUniFilter & AdaptKry [Hunag et al | ICML & WWV'24]

Application: Image Filtering

Graph Simplification

Future Direction

Spectral Embedding: Schema



Spectral Analysis

- Comprehensive graph information
- Explainable & controllable
- Separated graph computation

Spectral Embedding: Schema



Spectral Embedding

Channel Combination



How to select & combine useful channels?









Neural

Network



Spectral Embeddings

Channel Combination: LD² [NeurIPS'23]

• Heterophily: 1-hop message passing fails







Spectral Embedding

Channel Combination: LD² [NeurIPS'23]

• Heterophily: specifically focus 2-hop connections



Constant Feature:

Inverse Feature:

$$\boldsymbol{P}_{A} = \phi\left(\boldsymbol{A}^{2L} \cdot \boldsymbol{\mathcal{N}}\right)$$

$$oldsymbol{P}_{X,H} = \sum_{l=1}^{L} oldsymbol{ ilde{L}}^l \cdot oldsymbol{X}$$

 $oldsymbol{P}_{X,L2} = \sum_{l=1} oldsymbol{\bar{A}}^{2l} \cdot oldsymbol{X}$





N Liao et al, " LD²: Scalable Heterophilous Graph Neural Network with Decoupled Embedding", *NeurIPS'23*.

Channel Combination: LD² [NeurIPS'23]



Spectral Embedding

Spectral Embedding

Channel Combination: LD² [NeurIPS'23]



Decoupled spectral embedding



Channel Combination: S3GCL [ICML'24]

• Contrastive learning: difference between low and high frequency signals

Spectral Embedding



Channel Combination: TER [KDD'23]

• Representation based on edge-wise proximity:



• Eigen-decomposition:

$$\mathbf{Z} = \mathbf{U} \sqrt{\sum_{t=0}^{\infty} w(t)} \mathbf{\Lambda}^{t}$$

$$\boldsymbol{P} = \boldsymbol{U}\boldsymbol{\Lambda}\boldsymbol{U}^{\top}$$
$$\boldsymbol{\Pi} = \boldsymbol{U}\cdot\left(\sum_{t=0}^{\infty}w(t)\boldsymbol{\Lambda}^{t}\right)\cdot\boldsymbol{U}^{\top}$$

Spectral Embedding

Spectral Embedding

Adaptive Selection



Spectral Embedding

Adaptive Selection



How to filter based on learning feedback?



Spectral Embedding Adaptive Selection: UniFilter [ICML'24]



Spatial computation



K Huang et al, " How Universal Polynomial Bases Enhance Spectral Graph Neural Networks: Heterophily, Over-smoothing, and Over-squashing", ICML'24.



Adaptive Selection: AdaptKry [WWW'24]

• Understanding graph propagation as heat diffusion:

Laplacian

$$A = I - L$$

$$H_{t+1} = (I - L)H_t$$

$$\frac{\mathrm{d}H_t}{\mathrm{d}t} = -LH_t$$

$$\tau \to 0$$

$$H_{t+\tau} = \lim_{\tau \to 0^+} H_t - \tau \mathcal{L}H_t$$

$$= \lim_{\tau \to 0^+} H_t (I - \tau \mathcal{L})$$

$$= \lim_{\tau \to 0^+} H_t ((1 - \tau)I + \tau D^{-\frac{1}{2}}AD^{-\frac{1}{2}})$$

• Krylov Subspace: fixed propagation

$$\mathcal{K}_K(A,\mathbf{v}) = \{\mathbf{v}, A\mathbf{v}, A^2\mathbf{v}, \dots, A^{K-1}\mathbf{v}\}.$$



Spectral Embedding

Adaptive Selection: Adapt $|_{H_{t+\tau}} = \lim_{\tau \to 0^+} H_t - \tau \mathcal{L} H_t$ $= \lim_{t \to \infty} \mathbf{H}_t (\mathbf{I} - \tau \mathcal{L})$

• Renormalize propagation:

$$= \lim_{\tau \to 0^+} \mathbf{H}_t ((1 - \tau)\mathbf{I} + \tau \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}})$$

$$P = (1 - \tau)I + \tau D^{-1/2} A D^{-1/2}$$

$$P_{\tau} = (\tau D + (1 - \tau)I)^{-1/2} (\tau A + (1 - \tau)I) (\tau D + (1 - \tau)I)^{-1/2}$$

$$P_{\tau} = D_{\tau}^{-\frac{1}{2}} A_{\tau} D_{\tau}^{-\frac{1}{2}}$$

L

Adaptive Krylov Subspace: control graph message

$$\mathcal{K}_{K(\boldsymbol{P}_{\tau}, \mathbf{X})} = \{\mathbf{x}, \boldsymbol{P}_{\tau}\mathbf{x}, \boldsymbol{P}_{\tau}^{2}\mathbf{x}, \dots, \boldsymbol{P}_{\tau}^{K-1}\mathbf{x}\}$$

Spectral Embedding: Takeaways

Channel Combination

LD² [NeurIPS'23]

- 。 2-hop MP for global information
- . Link spatial & spectral computation

TER [KDD'23] S3GCL [ICML'24]

 Extending spectral embedding to edge-level task and contrastive learning

Adaptive Selection

UniFilter [ICML'24]

- . Link spatial & spectral computation
- Additional bias for hop-wise heterophily

AdaptKry [WWW'24]

- MP as heat diffusion
- 。 Adjust graph normalization

Spectral Embedding: Applications

Image Filtering

• Image as graph: better depict unstructured relationships, e.g. contextual connection of segmentations



Spectral Embedding: Evaluation

Graph Propagation as Spectral Filter

- Global view is important for large-scale prediction
- Simple, invariable filters
 can also capture global view

N Liao et al, "A Comprehensive Benchmark on Spectral GNNs: The Impact on Efficiency, Memory, and Effectiveness", SIGMOD'26.



Spectral Embedding: Evaluation

Graph Propagation as Spectral Filter

- Comprehensive evaluation on: efficiency, memory footprint, propagation hops of graph filters, ...
- Insights on finding fast and effective graph filters



Spectral Embedding: Evaluation

Weight Transformation as Spectral Filter

• Weight transformation also affects spectral filtering



Advances in Designing Scalable GNNs

Graph Simplification
Node-wise Similarity

Spectral Embedding

20 min 11:45 – 12:05 **Graph Simplification**

Ningyi Liao

Layer-wise NIGCN [K Huang et al | WWW'23]
Pair-wise ATP [X Li et al | WWW'24]
Entry-wise Unifiews [N Liao et al | ICML'25]

Future Direction

Graph Simplification: Strategies

Raw Graph





Sampling

- Different graphs through training iterations
- Repetitive sampling overhead



Sparsification

- Static graphs structure during learning
- Dedicated algorithms

Graph Simplification

Layer-wise Sparsification: NIGCN [WWW'23]

General Heat Diffusion

$$U(\omega,\rho,\ell) = \frac{\omega^\ell}{(\ell\,!)^\rho\cdot C}$$

- Coverage: smooth (GCN), PPR (GBP), HKPR (AGP)
- Applicable to decoupled architecture



Graph Simplification

Layer-wise Sparsification: NIGCN [WWW'23]





- Bounds computation complexity & approximation precision
- Nested in graph propagation

 $l_u = 2$

X Li et al, "Rethinking Node-wise Propagation for Large-scale Graph Learning", WWW'24.

Pair-wise Sparsification: ATP [WWW'24]

Conventional Graph Propagation

 $\Pi = \sum_{i=0}^{l} w_i \cdot \left(\hat{\mathbf{D}}^{r-1} \hat{\mathbf{A}} \hat{\mathbf{D}}^{-r} \right)^i,$

- Unified normalization for all nodes
- High-degree nodes produce biased accuracy



Graph Simplification

Graph Simplification

Pair-wise Sparsification: ATP [WWW'24]

Node-wise Propagation

$$\tilde{\Pi} = \sum_{i=0}^{l} w_i \cdot \left(\left[\hat{\mathbf{D}} \right]^{\tilde{\mathbf{R}}-1} \left[\hat{\mathbf{A}} \right] \left[\hat{\mathbf{D}} \right]^{-\tilde{\mathbf{R}}} \right)^i$$

$$\mathbf{R}_0 = \operatorname{Diag}\left(\frac{[d]_1}{n-1}, \dots, \frac{[d]_n}{n-1}\right).$$

- Specific normalization for each node
- Criteria based on degree, eigenvector, and connectivity



Pair-wise Sparsification: ATP [WWW'24]

Pair-wise Propagation Distance

$$\left| \left(\mathbf{P}^k e_i \right)_j - \tilde{\pi}_j \right| \le \sqrt{\frac{\tilde{d}_j}{\tilde{d}_i}} \lambda_2^k,$$

- Precompute the converged value
- Threshold propagations closed to converged values



Graph Simplification

X Li et al, "Rethinking Node-wise Propagation for Large-scale Graph Learning", WWW'24.

Sparsification: Unifews [ICML'25]

Graph Simplification

Neighborhood



Entry-wise Sparsification





C

Sparsification: Unifews [ICML'25]

Sparsification ⇔ Approximate Graph Smoothing

GNN Graph Smoothing:

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

Embedding Attribute Graph

Spectral Approximation:

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

Approx Bound

Error Guarantee:
 (general)
 (scale-free graph)

$$\Rightarrow$$
 $q_a \delta_a \leq \epsilon || p ||$
 \Rightarrow
 $\epsilon = O\left(\eta_a(1 - \eta_a)^{\frac{1}{1 - \alpha}}\right)$

 Sparsity
 Prune Threshold
 Graph Distribution



Sparsification: Unifews [ICML'25]

Layer-progressive Sparsification

Graph Sparsification:

 $\|\boldsymbol{p}_{(l+1)} - \hat{\boldsymbol{p}}_{(l+1)}\| \leq O(\epsilon) \cdot \|\boldsymbol{p}_{(l)}\|$ + $bc\epsilon(1+O(\epsilon)) \| \boldsymbol{p}_{(l)} \|$ Embedding / Approx Bound

w/ Weight Sparsification:

$$\begin{aligned} \|\hat{H}_{(l+1)} - H_{(l+1)}\|_{F} &\leq \ell_{\sigma}bc\epsilon \|H_{(l)}\|_{F} \|W\|_{F} \\ \text{Representation} &+ \ell_{\sigma}q_{w}\delta_{w} \\ & &$$

N Liao et al, " Unifews: You Need Fewer Operations for Efficient Graph Neural Networks", ICML'25.

Graph Sparsification G $\hat{\mathcal{G}}_{(1)}$ $\mathcal{G}_{(0)}$ Graph Representation Propagate \mathbf{r} $\mathbf{\mathbf{v}}\hat{\mathbf{T}}_{(1)}$ $\hat{H}_{(1)}$ $\hat{H}_{(2)}$ $\hat{H}_{(0)}$ Transform Attribute $(\delta_w$ Weight

 $\hat{W}_{(1)}$

X

Weight Sparsification

Graph Simplification



 $\hat{W}_{(2)}$

Graph Simplification: Takeaways

Layer-wise	Pair-wise	Entry-wise
NIGCN [WWW'23]	AIP [WWW'24]	Unitews [ICML'25]

Advances in Designing Scalable GNNs

Conclusion & Future Directions

Outline

Node-wise Similarity

Spectral Embedding

Graph Simplification

25 min 12:05 – 12:30 Future Directions Reynold Cheng

SummaryLarge ModelRAG | Graph TransformerData VariantData Deficiency | Dynamic DataSystem Co-designDevice-specific | Distributed Training

Conclusion

Summary: Challenges & Solutions



Reighbor Explosion



Node-wise Similarity

- PPR: faster message passing
- SimRank: global similarity
- Hub Labeling: core hub nodes

Limited Memory







Subgraph Extraction

- Partition: reduce transmission
- Generation: tailored message passing
- Slicing: disk lookup

Spectral Embedding

- Combination: local & global channels
- Selection: adaptive filtering

Graph + Large Models

Message-passing

- Mature & diverse options
- Tailored for graph data
- Scalable solutions

Transformer-based

- Powerful & attentive
 - Stackable to deep layers
 - Integrate with existing LMs



Model

LM: Graph for Retrieval-Augmented Generation

Standard RAG



LM: Graph for Retrieval-Augmented Generation

Opportunities for Graph + RAG



LM: Graph for Retrieval-Augmented Generation

Example: PageRank in RAG



N Alonso et al, "The Mixture-of-PageRanks Retriever for Long-Context Pre-Processing", 2024. https://www.zyphra.com/post/

LM: Scalable Graph Transformer Output Task $^{(+)}$ Raw Scale Chat Pointwise Pointwise Graph Feedforward Feedforward Pointwise Feedforward Layer Norm Layer Norm γ_2,β_2 Scale, Shift + Layer Norm Recom Multi-Head Multi-Head mend Self-Attention α_1 Layer Norm Scale Layer Norm (+)Multi-Head Self-Attention Multi-Head γ_1, β_1 \mathcal{O} Analyze Scale, Shift Self-Attention Concatenate on Sequence . . . Dimension Layer Norm MLP Laver Norm 888 886 Input Tokens Conditioning Input Tokens Input Tokens Conditioning . . . **Transformer Model** Graph Tokens Sequential Input



Efficient sequential graph representation?



Efficient graph Transformer architecture?

Graph Variants: Data Deficiency

Pretrain

Contrastive





Graph Variants: Data Deficiency

Example: Knowledge Sharing in Recommendation

- Challenge: Data Deficiency
 - Overlap Sparsity: few overlap data between domains
 - Data Sparsity: limited data available in some domains
- Solution: Cross-domain learning by knowledge alignment



J Liu et al. "Cross-domain knowledge graph chiasmal embedding for multi-domain item-item recommendation". *TKDE*'22 W Ning et al. "Multi-domain Recommendation with Embedding Disentangling and Domain Alignment". *CIKM*'23.

Graph Variants: Data Deficiency

Example: Knowledge Sharing in Recommendation

- Challenge: Model Scalability
 - Inter-domain: effectively learn domainspecific knowledge
 - Intra-domain: efficiently share knowledge without learning from scratch
- Solution: Model disentangling
 - Domain Model
 - Shared Model



Graph Variants: Dynamic Data

Systems Co-design: Device-Specific

Graph Storage

Graph Computation

On-disk

- Large storage
- Transmission overhead
- Update overhead

GPU-oriented

- Limited memory
- Batch processing
- Predefined operators

In-memory

- Random access
- Device loading
- 。Cache & buffer

CPU-oriented

- Sequential/parallel
- Dedicated algorithms
- Sparse data structures

Systems Co-design: Distributed Training

- Neighbors affected by partitioning
- More communication for
 - accessing neighbors

Multi-scale

- Local subgraph structure affected by partitioning
- Communication for global info

THANK YOU

2025 International Conference on Management of Data Berlin, Germany – June 22, 2025