



# Self-supervised Learning and Pre-training on Graphs (**GNNs**)

Yukuo Cen, Yuxiao Dong, Jie Tang

Knowledge Engineering Group (KEG)

Department of Computer Science and Technology

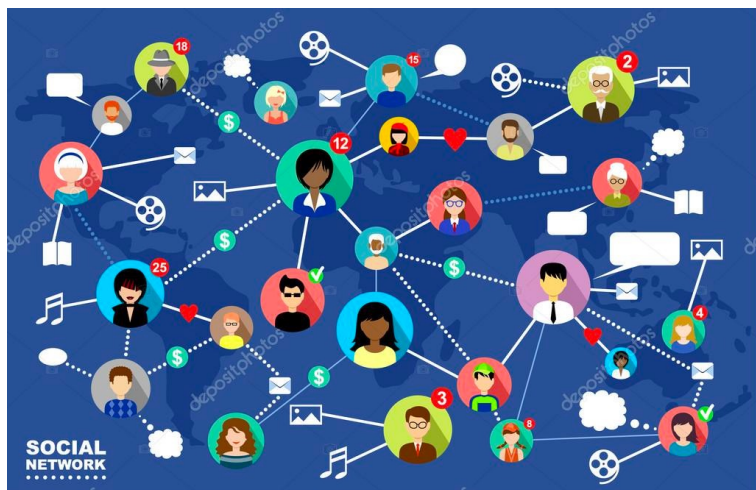
Tsinghua University



[Download the slides here](#)

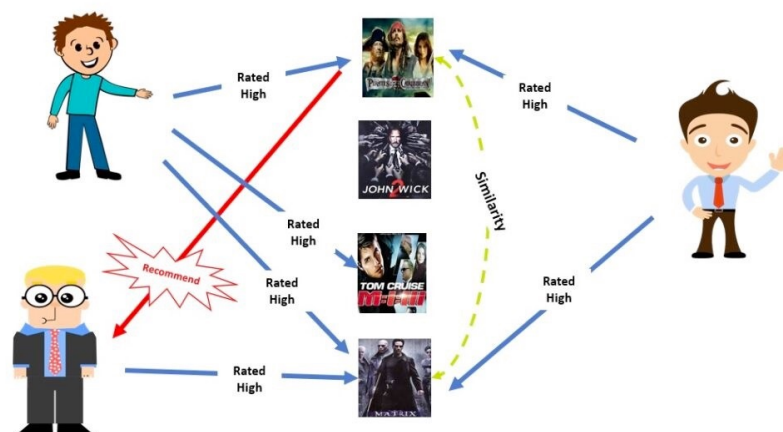
# Graph

- **Graph data** exists everywhere



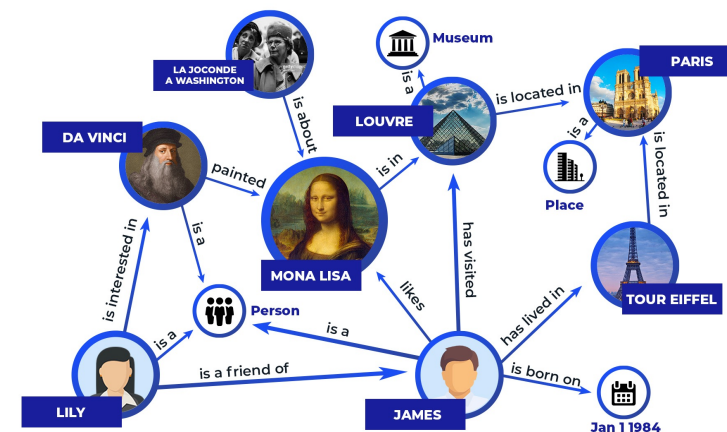
Social Network

- WeChat: 1.2 billion users  
61 billion links



Recommender System

- Alibaba: 2.3 billion trans. on 11/11



Knowledge Graph

- Wikidata: >1.4 billion triples

*“The number of **graph neural network** papers in this journal has grown as the field matures. We take a closer look at some of the **scientific applications**.”*

# Machine Learning on Graphs

- ML tasks on Graphs:
  - Node classification
    - Predict a type of a given node
  - Link prediction
    - Predict whether two nodes are linked
  - Graph classification
    - Predict the properties of molecules
  - Community detection
    - Identify densely linked clusters of nodes

Learning on Graphs with Graph Neural Networks (GNNs)

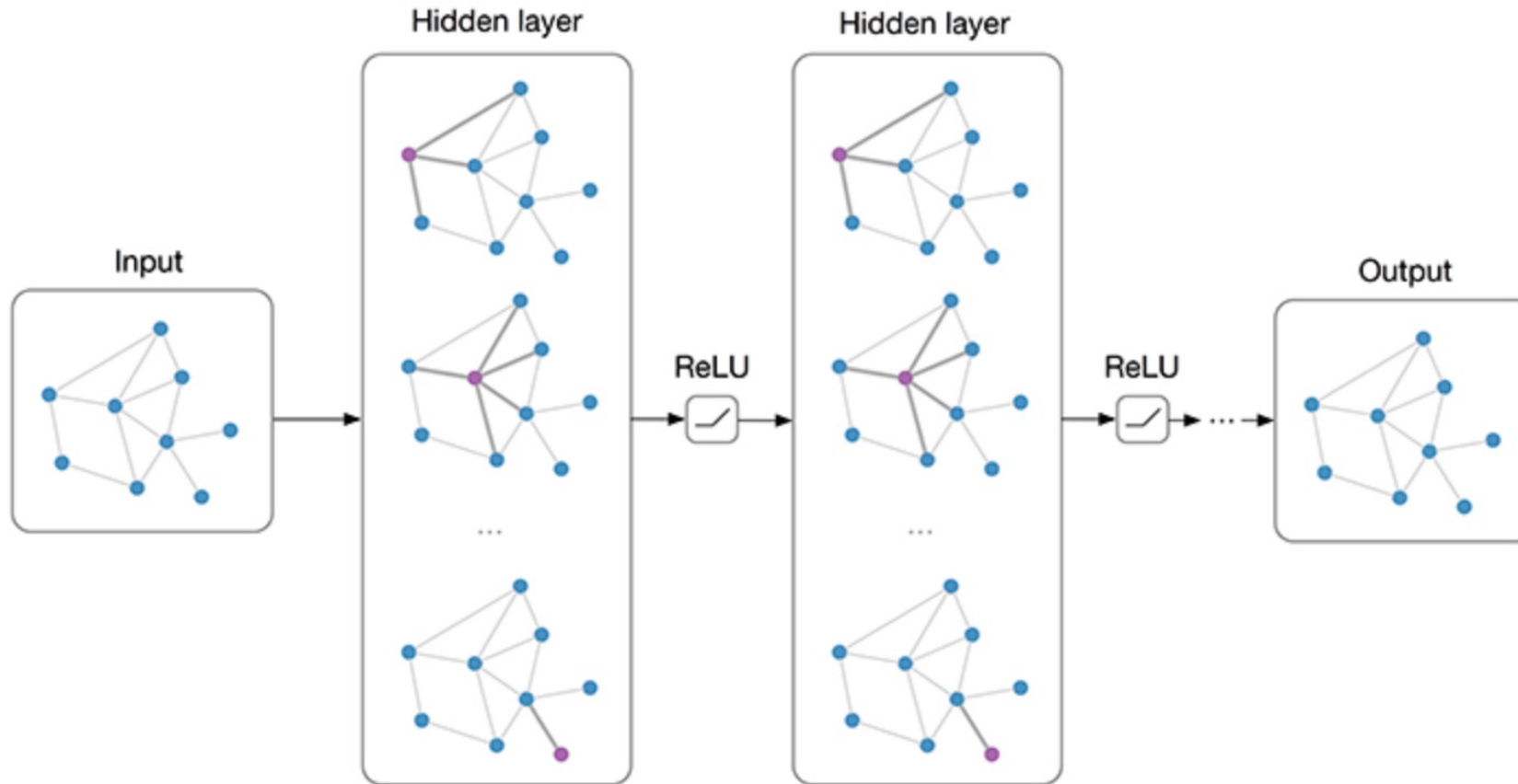
- A question: *Are you using GNNs?*



# Graph Neural Networks

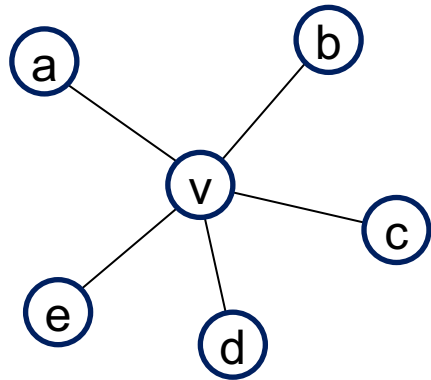
- Layer-wise propagation:

$$f(H^{(l)}, A) = \sigma(AH^{(l)}W^{(l)})$$



Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In ICLR '17.

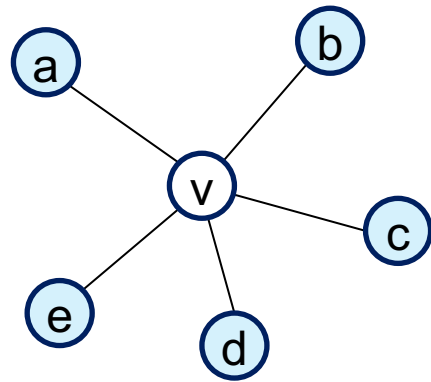
# Graph Neural Networks



$$h_v = f(h_a, h_b, h_c, h_d, h_e)$$

- **Neighborhood Aggregation:**
  - Aggregate neighbor information and pass into a neural network
  - It can be viewed as a center-surround filter in CNN---graph convolutions!

# GCN: Graph Convolutional Networks



parameters in layer  $k$

Non-linear activation function (e.g., ReLU)

$$\mathbf{h}_v^k = \sigma \left( \mathbf{W}_k \sum_{u \in \mathbf{N}(v) \cup v} \frac{h_u^{k-1}}{\sqrt{|N(u)||N(v)|}} \right)$$

node  $v$ 's embedding at layer  $k$

the neighbors of node  $v$

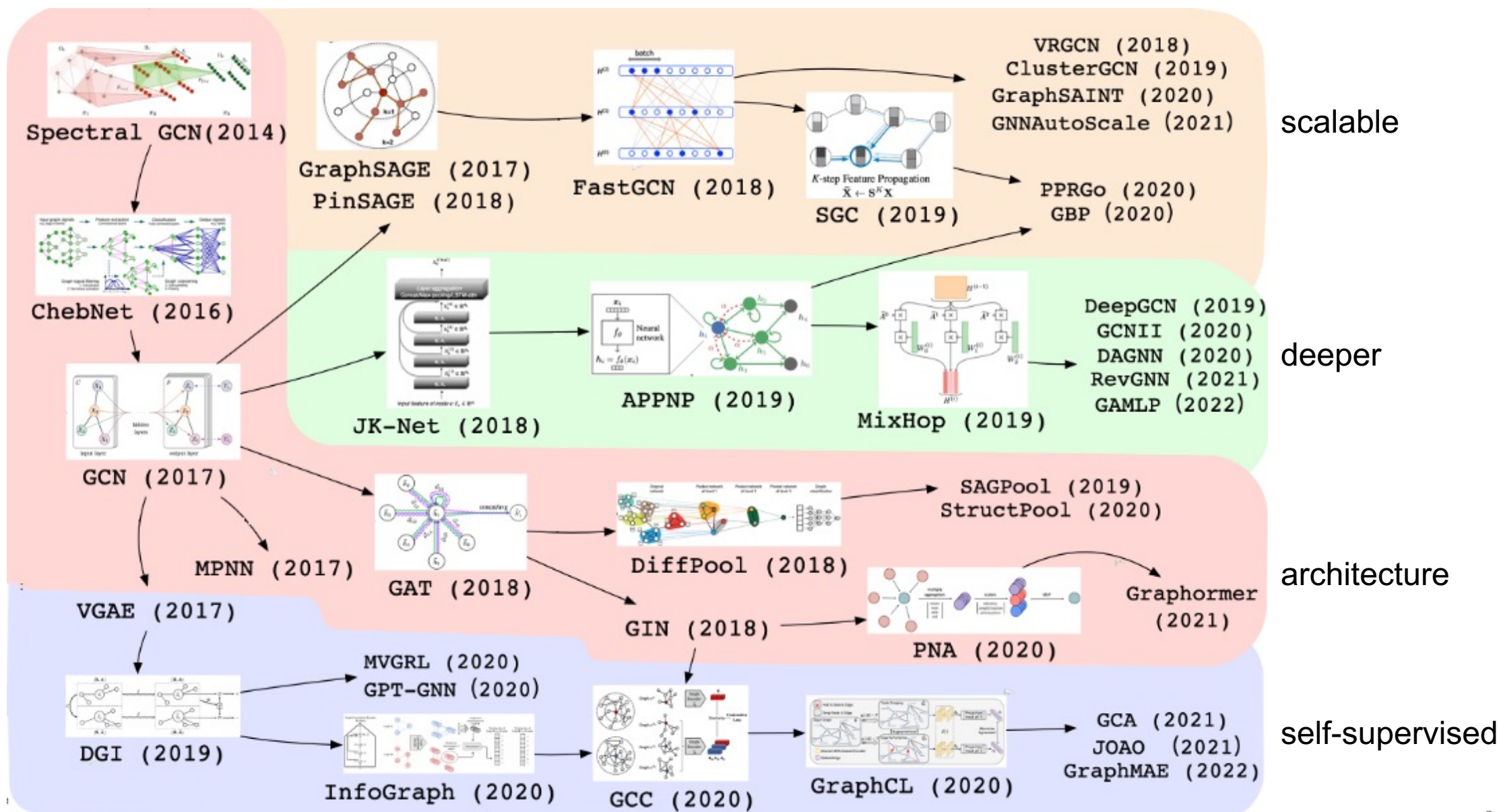
# GCN Performance

- 2-layer GCN:  $\mathbf{Z} = \text{softmax}(\tilde{\mathbf{A}} \sigma(\tilde{\mathbf{A}} \mathbf{X} \mathbf{W}_0) \mathbf{W}_1)$

Dataset	Type	Nodes	Edges	Classes	Features	Label rate
Citeseer	Citation network	3,327	4,732	6	3,703	0.036
Cora	Citation network	2,708	5,429	7	1,433	0.052
Pubmed	Citation network	19,717	44,338	3	500	0.003
NELL	Knowledge graph	65,755	266,144	210	5,414	0.001

Method	Citeseer	Cora	Pubmed	NELL
ManiReg [3]	60.1	59.5	70.7	21.8
SemiEmb [28]	59.6	59.0	71.1	26.7
LP [32]	45.3	68.0	63.0	26.5
DeepWalk [22]	43.2	67.2	65.3	58.1
ICA [18]	69.1	75.1	73.9	23.1
Planetoid* [29]	64.7 (26s)	75.7 (13s)	77.2 (25s)	61.9 (185s)
<b>GCN (this paper)</b>	<b>70.3 (7s)</b>	<b>81.5 (4s)</b>	<b>79.0 (38s)</b>	<b>66.0 (48s)</b>

# GNN History



# Do we really make big progress?

- Using “heterogeneous graph neural networks (HGNN)” as an example
- **Unrobust** results with **biased** setting on **small** data

	HAN [36]		GTN [43]			RSHN [45]			HetGNN [44]				MAGNN [12]	
Dataset	ACM		DBLP	ACM	IMDB	AIFB	MUTAG	BGS	MC (10%)		MC (30%)		DBLP	
Metric	Macro-F1	Micro-F1	Macro-F1	Macro-F1	Macro-F1	Accuracy	Accuracy	Accuracy	Macro-F1	Micro-F1	Macro-F1	Micro-F1	Macro-F1	Micro-F1
model*	91.89	91.85	94.18	92.68	60.92	97.22	82.35	93.10	97.8	97.9	98.1	98.2	93.13	93.61
GCN*	89.31	89.45	87.30	91.60	56.89	-	-	-	-	-	-	-	88.00	88.51
GAT*	90.55	90.55	93.71	92.33	58.14	91.67	72.06	66.32	96.2	96.3	96.5	96.5	91.05	91.61
model	We tested 12 HGNN algorithms													
GCN														
GAT														

**\* With a fairly proper setting, the results are even reversed!**

# Challenges

- Challenge 1: Robustness
- Challenge 2: Unlabeled data
- Challenge 3: Easy-to-use Tool

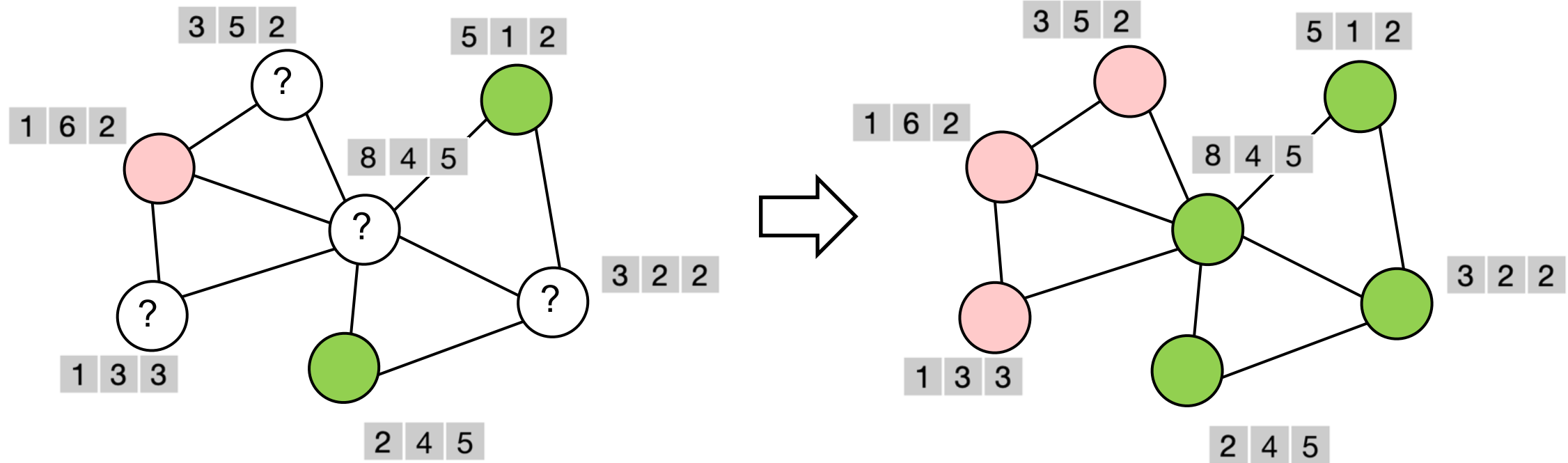


# Overview

- Semi-supervised Learning on Graphs:
  - Training data: a small portion of labeled data + lots of unlabeled data
  - **Robustness**: consistency regularization for predictions of different views
- CogDL: A Comprehensive Library for Graphs (**Easy-to-use**)
- Contrastive Self-supervised Learning on Graphs:
  - Training data: all data is **unlabeled**
  - Contrasts the views generated from different augmentations
- Generative Self-supervised Learning on Graphs:
  - Training data: all data is **unlabeled**
  - Reconstruction of the input graph (graph structure, node features)

# Semi-supervised Learning on Graphs

# Semi-Supervised Learning on Graphs

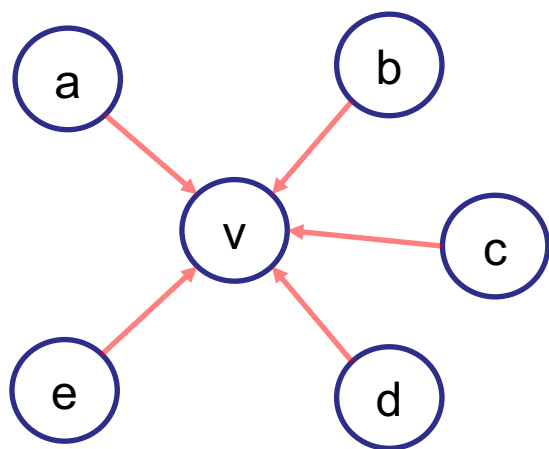


**Input:** a partially labeled & attributed graph

**Output:** infer the labels of unlabeled nodes

# Graph Neural Networks (GNN)

Message Passing:



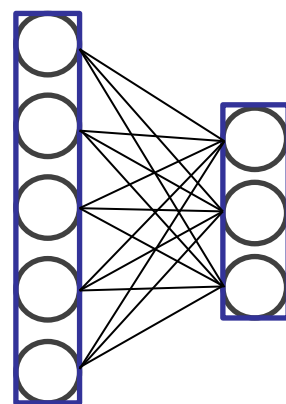
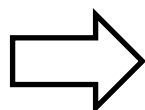
Feature propagation

Representation vector of the  $l + 1$ -th layer

Activation function (e.g. ReLU)

GCN:  $\mathbf{H}^{l+1} = \sigma(\hat{\mathbf{A}}\mathbf{H}^{(l)}\mathbf{W}^{(l)})$

Normalized adjacency matrix



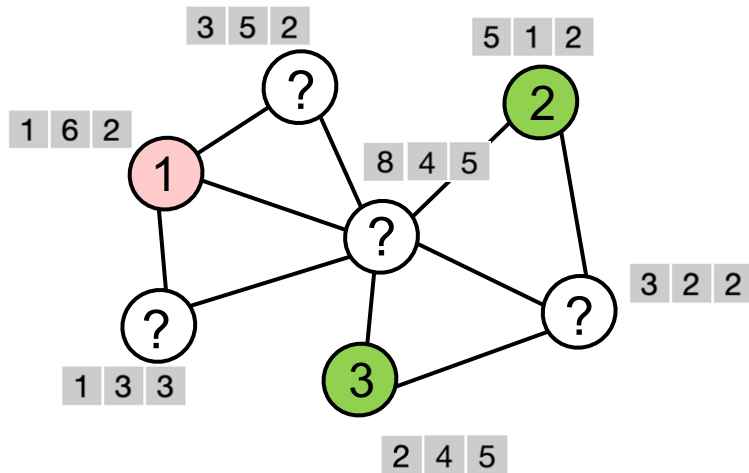
Non-linear transformation

$$\mathbf{H}^{l+1} = \sigma \left( \mathbf{W}^l \sum_{u \in N(v) \cup v} \frac{\mathbf{H}_u^l}{\sqrt{|N(u)| |N(v)|}} \right)$$

# Over-fitting problem of GNNs

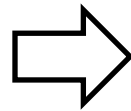
1. In GNNs, feature propagation is coupled with non-linear transformation. Increasing layers will introduce unnecessary parameters.
2. GNNs only adopt the supervised cross-entropy loss to guide the model training.

Training process:

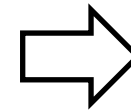


Feature propagation is coupled with  
Non-linear transformation

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



GNN



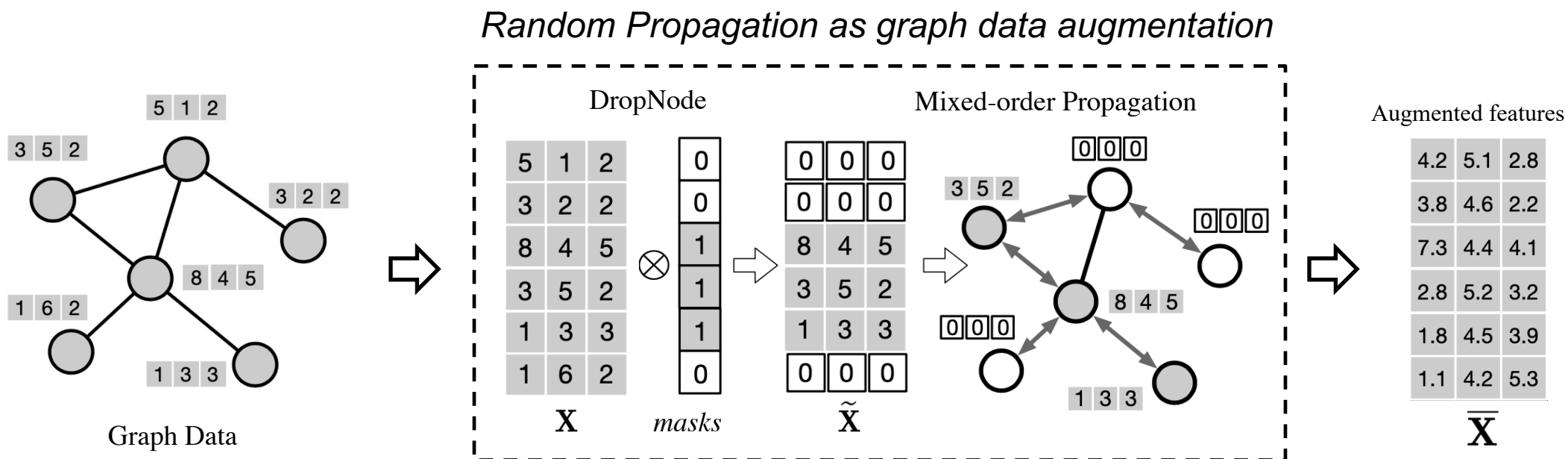
Loss function:

$$\mathbf{y}_1^T \log(\hat{\mathbf{y}}_1) + \mathbf{y}_2^T \log(\hat{\mathbf{y}}_2) + \mathbf{y}_3^T \log(\hat{\mathbf{y}}_3)$$

Cannot fully leverage the  
unlabeled nodes

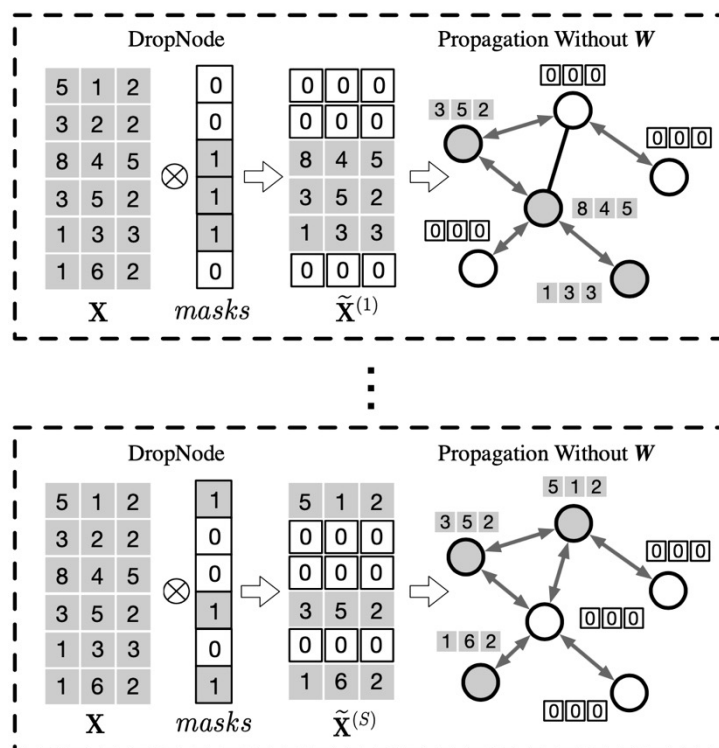
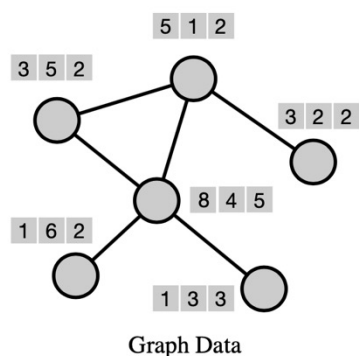
# Graph Random Neural Network (GRAND)

- Random Propagation (DropNode + Propagation):
  - Decouple the feature propagation from non-linear feature transformation.
  - Propagate feature with a mixed-order adjacency matrix:  $\Pi = \sum_{n=0}^N \frac{1}{N+1} \hat{\mathbf{A}}^n$
  - Use DropNode before feature propagation to randomly aggregate neighbors' features



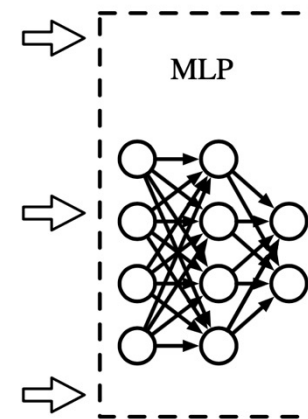
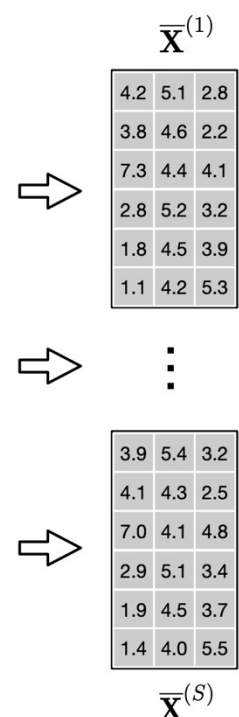
# Graph Random Neural Network (GRAND) (cont.)

## $M$ Augmentations

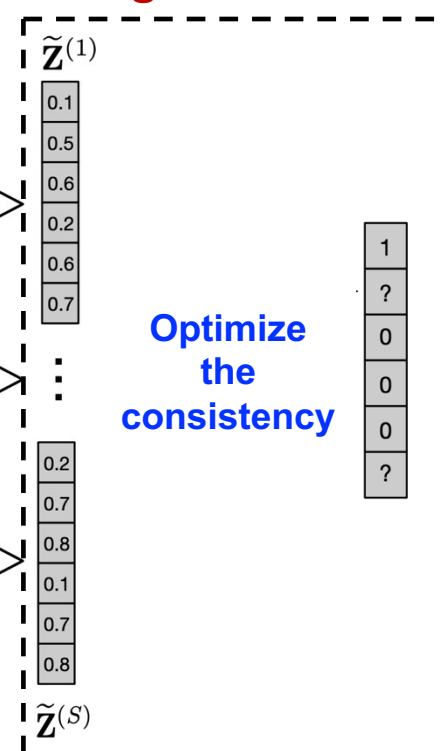


Random Propagation

Augmented features  $\bar{\mathbf{X}}$



## Consistency Regularization

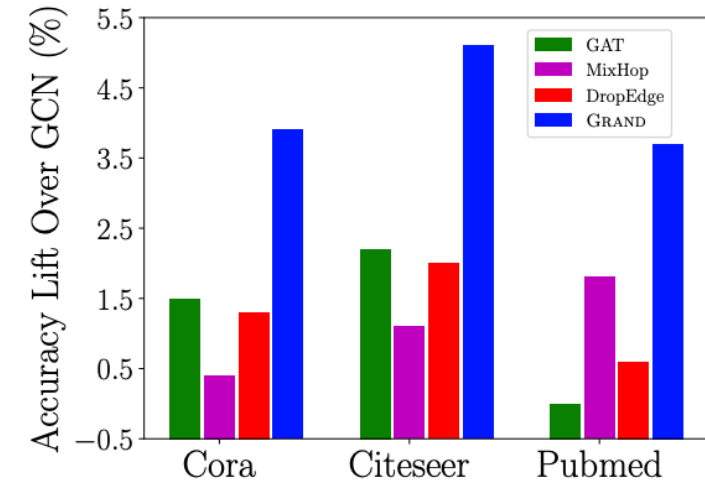


Fully leverage the unlabeled data



# Performance of GRAND

	Method	Cora	Citeseer	Pubmed
GCNs	GCN [19]	81.5	70.3	79.0
	GAT [32]	83.0 $\pm$ 0.7	72.5 $\pm$ 0.7	79.0 $\pm$ 0.3
	APPNP [20]	83.8 $\pm$ 0.3	71.6 $\pm$ 0.5	79.7 $\pm$ 0.3
	Graph U-Net [11]	84.4 $\pm$ 0.6	73.2 $\pm$ 0.5	79.6 $\pm$ 0.2
	SGC [36]	81.0 $\pm$ 0.0	71.9 $\pm$ 0.1	78.9 $\pm$ 0.0
	MixHop [1]	81.9 $\pm$ 0.4	71.4 $\pm$ 0.8	80.8 $\pm$ 0.6
	GMNN [28]	83.7	72.9	81.8
	GraphNAS [12]	84.2 $\pm$ 1.0	73.1 $\pm$ 0.9	79.6 $\pm$ 0.4
Sampling GCNs	GraphSAGE [16]	78.9 $\pm$ 0.8	67.4 $\pm$ 0.7	77.8 $\pm$ 0.6
	FastGCN [7]	81.4 $\pm$ 0.5	68.8 $\pm$ 0.9	77.6 $\pm$ 0.5
Regularization GCNs	VBAT [10]	83.6 $\pm$ 0.5	74.0 $\pm$ 0.6	79.9 $\pm$ 0.4
	G <sup>3</sup> NN [24]	82.5 $\pm$ 0.2	74.4 $\pm$ 0.3	77.9 $\pm$ 0.4
	GraphMix [33]	83.9 $\pm$ 0.6	74.5 $\pm$ 0.6	81.0 $\pm$ 0.6
	DropEdge [29]	82.8	72.3	79.6
	GRAND	<b>85.4<math>\pm</math>0.4</b>	<b>75.4<math>\pm</math>0.4</b>	<b>82.7<math>\pm</math>0.6</b>



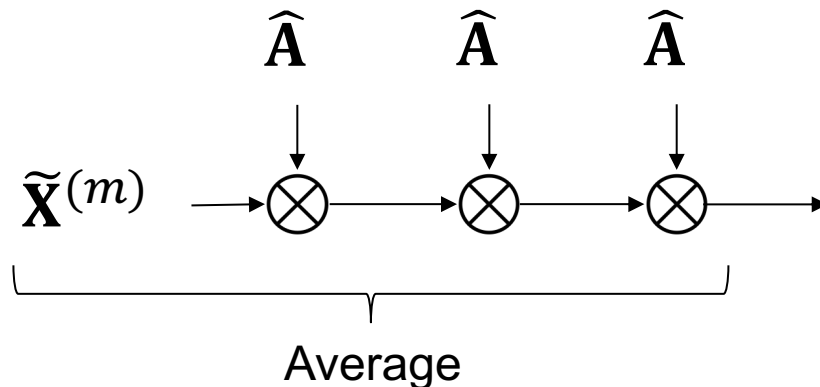
Instead of the marginal improvements by conventional GNN baselines over GCN, **GRAND** achieves ***much more significant performance lift in all three datasets!***

# Scalability limitation of GRAND

- Random Propagation in GRAND:

$$\bar{\mathbf{X}}^{(m)} = \Pi \tilde{\mathbf{X}}^{(m)}, \quad \Pi = \sum_{n=0}^N \frac{1}{N+1} \hat{\mathbf{A}}^n$$

- $\bar{\mathbf{X}}^{(m)}$  is calculated with power iteration:



## Weak Scalability:

- Time/Memory complexity:  $O(|E| + |V|)$ .
- Random propagation needs to be formed for multiple times at each epoch.

# GRAND+: General Idea

- Mini-batch Radom Propagation:
  - Select a batch of nodes at each training step, and generate augmented features by

$$\bar{\mathbf{X}}_s^{(m)} = \sum_{v \in \mathcal{N}_v^\pi} \mathbf{z}_v \cdot \Pi(s, v) \cdot \mathbf{X}_v, \quad \mathbf{z}_v \sim \text{Bernoulli}(1 - \delta)$$

DropNode mask

Non-zero elements in  $\Pi_s$

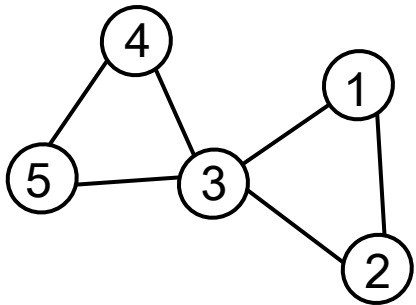
$$\Pi = \sum_{n=0}^N \frac{1}{N+1} \hat{\mathbf{A}}^n$$

How to efficiently calculate the row vector  $\Pi_s$  ?

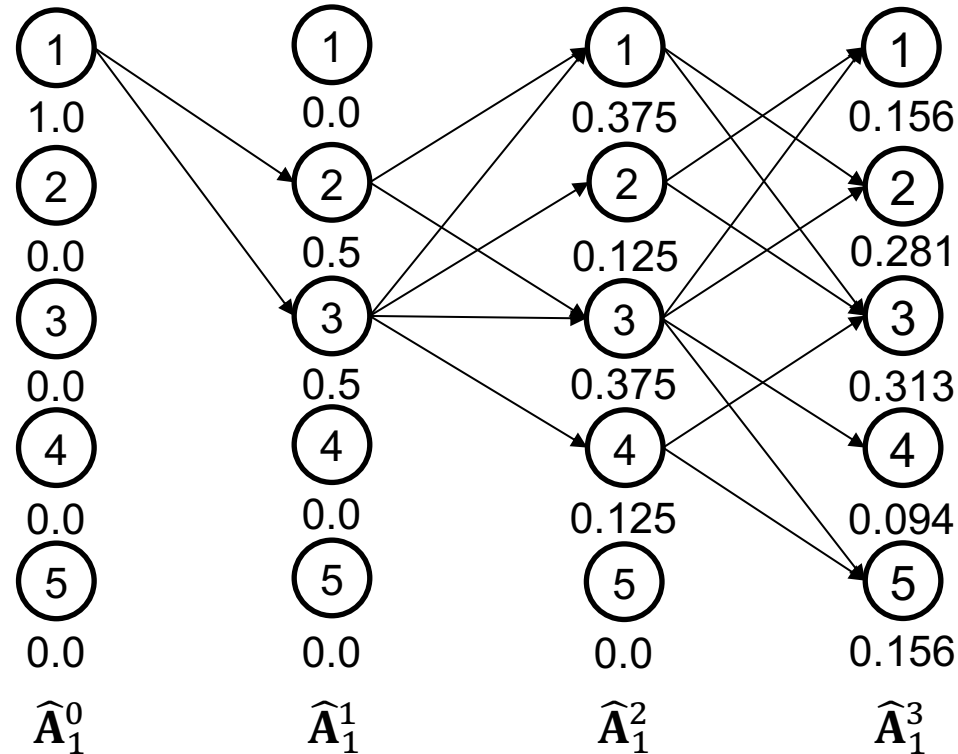
# GRAND+: Matrix approximation

$$\Pi = \frac{1}{N+1} \sum_{n=0}^N \hat{\mathbf{A}}^n$$

$\hat{\mathbf{A}} = \tilde{\mathbf{D}}^{-1} \tilde{\mathbf{A}}$  is random walk reverse transition matrix.  
 $\mathbf{P}(s, v)$  indicates the random walk probability from  $s$  to  $v$ .



$$\Pi_1 = \frac{1}{N+1} \sum_{n=0}^N \hat{\mathbf{A}}_1^n$$



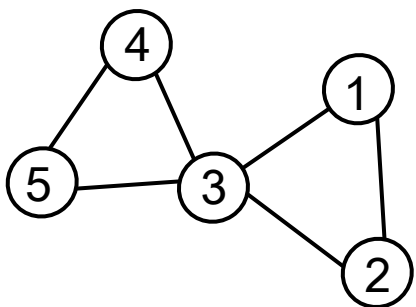
Random Walk Probability Diffusion

Complexity:  $O(|E|)$  ☹️

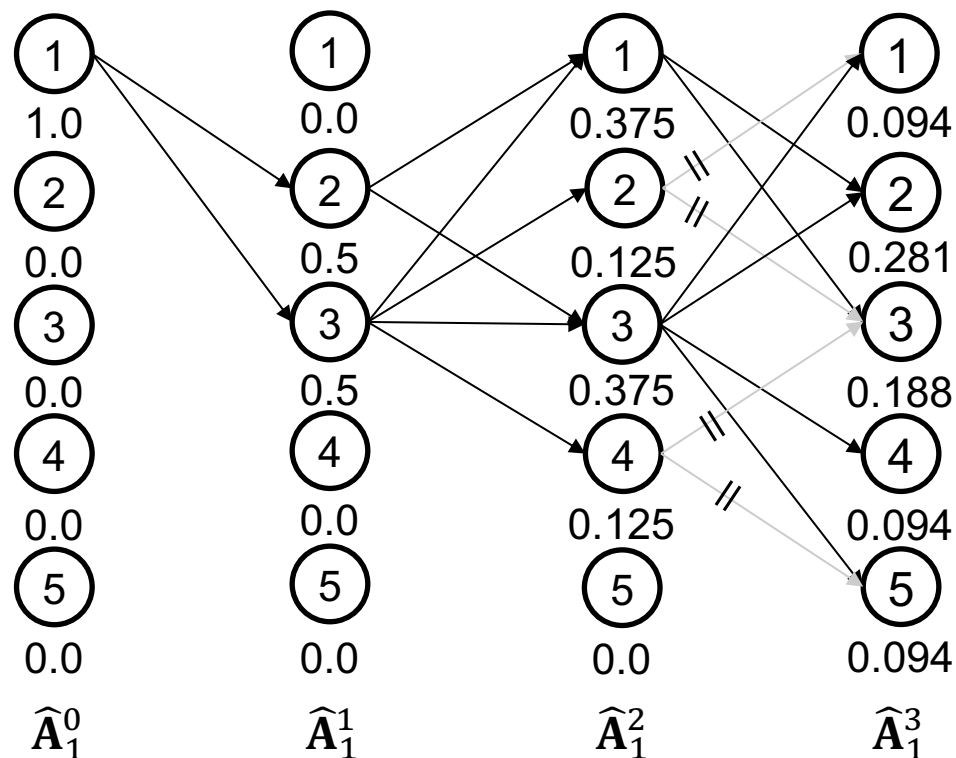
# GRAND+: Matrix approximation

$$\Pi = \frac{1}{N+1} \sum_{n=0}^N \hat{\mathbf{A}}^n$$

$\hat{\mathbf{A}} = \tilde{\mathbf{D}}^{-1} \tilde{\mathbf{A}}$  is random walk reverse transition matrix.  
 $\mathbf{P}(s, v)$  indicates the random walk probability from  $s$  to  $v$ .



$$\tilde{\Pi}_1 = \frac{1}{N+1} \sum_{n=0}^N \hat{\mathbf{A}}_1^n$$



Generalized Forward Push (GFPush)

$$r_{max} = 0.07$$

Error Bound:  $T \cdot r_{max}$

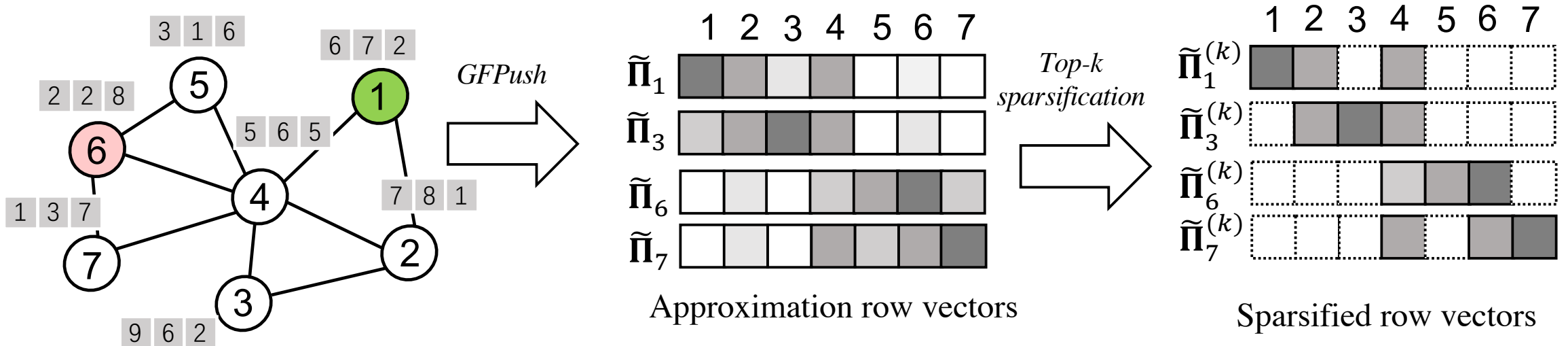
Time Complexity:  $T/r_{max}$

Memory Complexity:  $T/r_{max}$

# GRAND+: Matrix approximation

- Approximation method:
  - GFPush: Generate an error bounded approximation  $\tilde{\Pi}_s$  for  $\Pi_s$ .
  - Top- $k$  sparsification: Truncate  $\tilde{\Pi}_s$  to make it only contains top- $k$  elements.

**THEOREM 1.** *Algorithm 1 has  $O(N/r_{max})$  time complexity and  $O(N/r_{max})$  memory complexity, and returns  $\tilde{\Pi}_s$  as an approximation of  $\Pi_s$  with the  $L_1$  error bound:  $\|\Pi_s - \tilde{\Pi}_s\|_1 \leq N \cdot (2|E| + |V|) \cdot r_{max}$ .*



# GRAND+: Mini-batch Radom Propagation

- Mini-batch Radom Propagation with Approximation:

$$\bar{\mathbf{X}}_s^{(m)} = \sum_{v \in \mathcal{N}_v^{(k)}} \mathbf{z}_v \cdot \tilde{\Pi}^{(k)}(s, v) \cdot \mathbf{X}_v, \quad \mathbf{z}_v \sim \text{Bernoulli}(1 - \delta)$$

Non-zero elements in  $\tilde{\Pi}_v^{(k)}$

- Prediction:

$$\hat{\mathbf{Y}}^{(m)} = \text{MLP}(\bar{\mathbf{X}}_s^{(m)}, \Theta)$$

With batch size as  $b$ , the time complexity is  $O(b \cdot k)$ , which is independent of graph size

**Scalability:** Adopt GFPush to approximately calculate the propagation matrix, and adopt mini-batch method for model training



# GRAND+: Confidence-aware Consistency Regularization

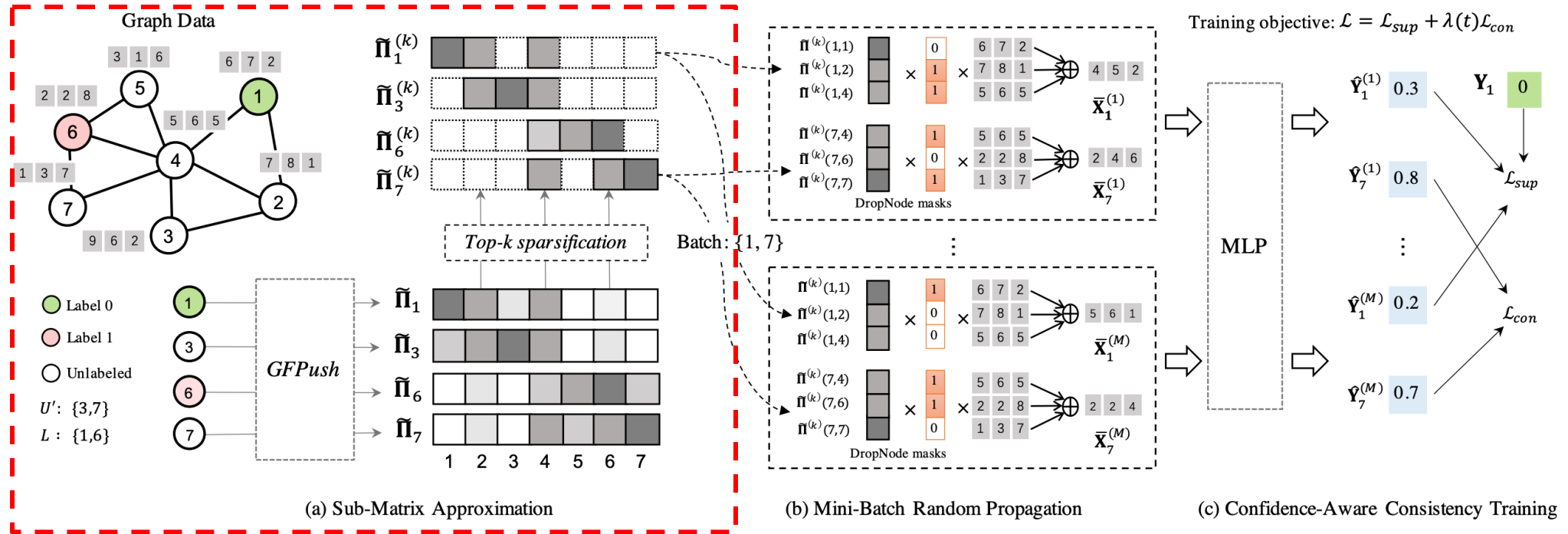
- Confidence-aware Consistency Loss:

$$\mathcal{L}_{con} = \frac{1}{b_u \cdot M} \sum_{s \in U_n} \mathbb{1}(\max(\bar{Y}_s) \geq \gamma) \sum_{m=1}^M \mathcal{D}(\tilde{Y}_s, \hat{Y}_s^{(m)}),$$

Confidence term: Filter out unlabeled nodes that have low confidence

**Effectiveness:** Further improving prediction performance

# GRAND+ Architecture



Parallelization by OpenMP

GRAND+: Better scalability & generalization capability

# Experiments

**Table 2: Classification Accuracy (%) on Benchmarks.**

Category	Method	Cora	Citeseer	Pubmed
Full-batch GNNs	GCN	$81.5 \pm 0.6$	$71.3 \pm 0.4$	$79.1 \pm 0.4$
	GAT	$83.0 \pm 0.7$	$72.5 \pm 0.7$	$79.0 \pm 0.3$
	APPNP	$84.1 \pm 0.3$	$71.6 \pm 0.5$	$79.7 \pm 0.3$
	GCNII	$85.5 \pm 0.5$	$73.4 \pm 0.6$	$80.3 \pm 0.4$
	GRAND	$85.4 \pm 0.4$	$75.4 \pm 0.4$	$82.7 \pm 0.6$
Scalable GNNs	FastGCN	$81.4 \pm 0.5$	$68.8 \pm 0.9$	$77.6 \pm 0.5$
	GraphSAINT	$81.3 \pm 0.4$	$70.5 \pm 0.4$	$78.2 \pm 0.8$
	SGC	$81.0 \pm 0.1$	$71.8 \pm 0.1$	$79.0 \pm 0.1$
	GBP	$83.9 \pm 0.7$	$72.9 \pm 0.5$	$80.6 \pm 0.4$
	PPRGo	$82.4 \pm 0.2$	$71.3 \pm 0.3$	$80.0 \pm 0.4$
Our Methods	GRAND+ (P)	<b><math>85.8 \pm 0.4</math></b>	<b><math>75.6 \pm 0.4</math></b>	$84.5 \pm 1.1$
	GRAND+ (A)	$85.5 \pm 0.4$	$75.5 \pm 0.4$	<b><math>85.0 \pm 0.6</math></b>
	GRAND+ (S)	$85.0 \pm 0.5$	$74.4 \pm 0.5$	$84.2 \pm 0.6$

Better generalization performance: Achieves **2.3%** improvements over GRAND on Pubmed.

# Experiments

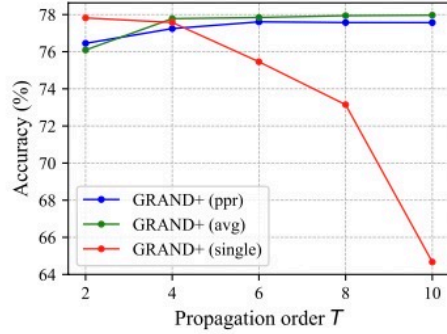
**Table 3: Accuracy (%) and Running Time (s) on Large Graphs.**

Method	AMiner-CS		Reddit		Amazon2M		MAG.	
	Acc	RT	Acc	RT	Acc	RT	Acc	RT
GRAND	53.1±1.1	750	OOM	–	OOM	–	OOM	–
FastGCN	48.9±1.6	69	89.6±0.6	158	72.9±1.0	239	64.3±5.6	4220
GraphSAINT	51.8±1.3	39	92.1±0.5	39	75.9±1.3	189	75.0±1.7	6009
SGC	50.2±1.2	9	92.5±0.2	31	74.9±0.5	69	–	>24h
GBP	52.7±1.7	21	88.7±1.1	370	70.1±0.9	280	–	>24h
PPRGo	51.2±1.4	11	91.3±0.2	233	67.6±0.5	160	72.9±1.1	434
GRAND+ (P)	53.9±1.8	17	93.3±0.2	183	75.6±0.7	188	77.6±1.2	653
GRAND+ (A)	54.2±1.7	14	<b>93.5±0.2</b>	174	75.9±0.7	136	<b>80.0±1.1</b>	737
GRAND+ (S)	<b>54.2±1.6</b>	10	92.8±0.2	62	<b>76.2±0.6</b>	80	77.8±0.9	483

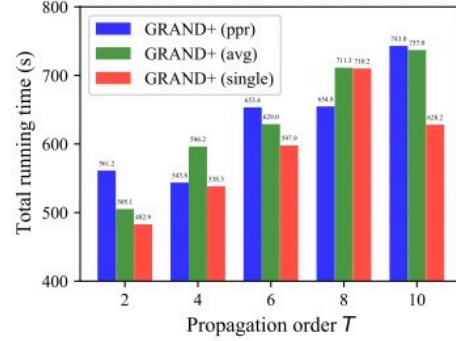
## Scalability:

- **40** times faster than GRAND on Aminer-CS.
- **8** times faster than FastGCN on MAG.
- **12** times faster than GraphSAINT on MAG.
- Achieves comparable running time and **4.9%** improvement than PPRGo on MAG.

# Parameter Analysis

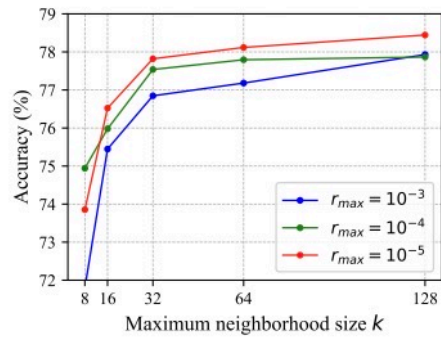


(a) Classification accuracy.

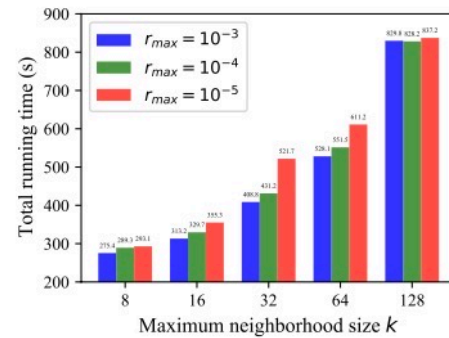


(b) Total running time.

**Figure 5: Effects of propagation order  $T$  on MAG-Scholar-C.**

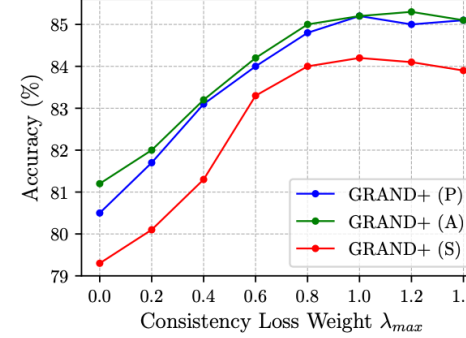


(a) Classification accuracy.

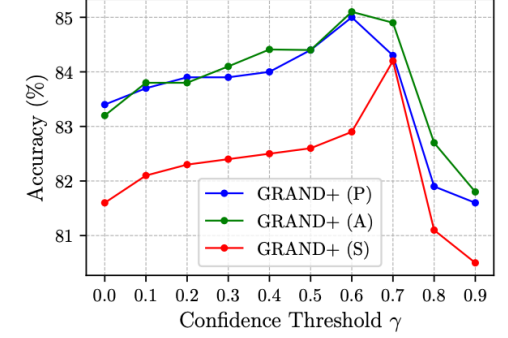


(b) Total running time.

**Figure 4: GRAND+ w.r.t.  $k$  and  $r_{max}$  on MAG-Scholar-C.**

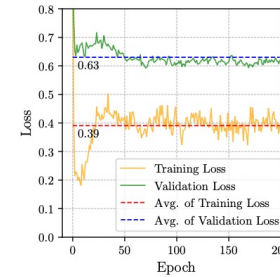


(a) Accuracy w.r.t.  $\lambda_{max}$ .

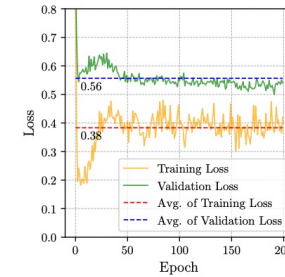


(b) Accuracy w.r.t.  $\gamma$

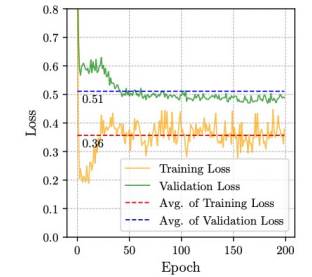
**Figure 2: Effects of  $\lambda_{max}$  and  $\gamma$  on Pubmed.**



(a)  $\lambda_{max} = 0.0$



(b)  $\lambda_{max} = 1.0, \gamma = 0.0$



(c)  $\lambda_{max} = 1.0, \gamma = 0.6$

**Figure 3: Training and Validation Losses on Pubmed.**

# CogDL: A Comprehensive Library for Graph Deep Learning



Access the code here

GitHub: <https://github.com/THUDM/cogdl> , 1400+ stars as of April 2023.

# CogDL – Overview

## Vision

CogDL aims at providing researchers and practitioners with easy-to-use APIs, reproducible results, and high efficiency for most graph tasks and applications.

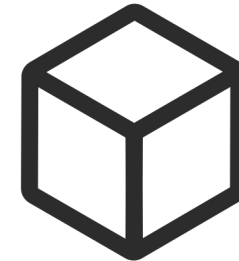
## Philosophy



Easy-to-use



Reproducibility



Efficiency



# CogDL - A Unified GNN Trainer

- Design a unified Trainer for GNN training



```
1 model = GCN(...)
2 data = CoraDataset()[0]
3 mw = ModelWrapper(model, ...)
4 dw = DataWrapper(data, ...)
5 trainer = Trainer(epochs=100)
6 result = trainer.run(mw, dw)
```

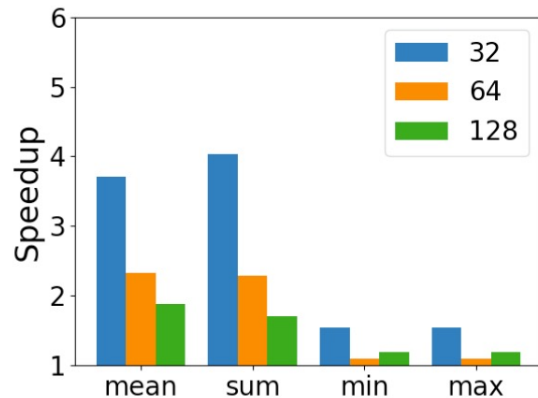
```
model = GCN(...)
data = Planetoid(name="Cora")[0]
optimizer = torch.optim.Adam(...)
for epoch in range(100):
    pred = model(data.x, data.edge_index)
    labels = data.y
    mask = data.train_mask
    loss = F.nll_loss(pred[mask],
                      labels[mask])
    optimizer.zero_grad()
    loss.backward()
    optimizer.step()
    val_acc = evaluate(model, data)
    if val_acc > best_val_acc:
        best_val_acc = val_acc
        best_model = deepcopy(model)
result = test(best_model, data)
```

```
model = GCN(...)
data = CoraGraphDataset()[0]
optimizer = torch.optim.Adam(...)
for epoch in range(100):
    pred = model(data)
    labels = data.ndata['label']
    mask = data.ndata['train_mask']
    loss = F.nll_loss(pred[mask],
                      labels[mask])
    optimizer.zero_grad()
    loss.backward()
    optimizer.step()
    val_acc = evaluate(model, data)
    if val_acc > best_val_acc:
        best_val_acc = val_acc
        best_model = deepcopy(model)
result = test(best_model, data)
```

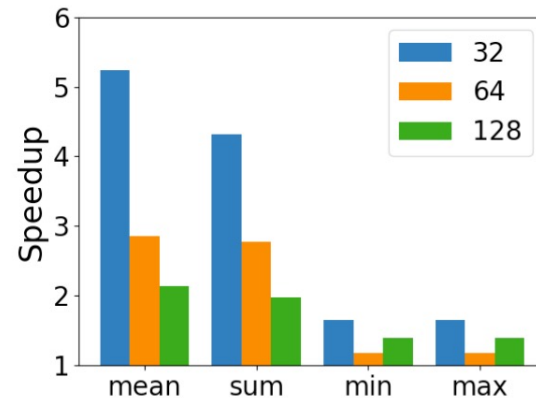
# CogDL – Operator Acceleration

- Acceleration for SpMM-like operators
  - Design a more balanced strategy for GPU parallel computation
  - **1.17x~5.24x** Speedup on Reddit/Yelp datasets compared to DGL

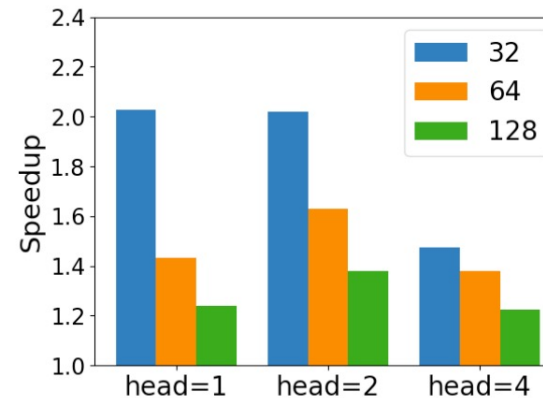
```
workload_balance = balancer(graph) # preprocessing
out_feat = GSpMM(workload_balance, graph, in_feat, reduce_op,
                  compute_op)
```



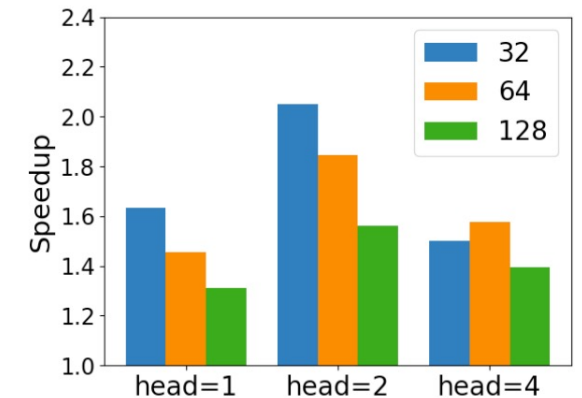
(a) GSpMM on Reddit



(b) GSpMM on Yelp



(c) multi-head SpMM on Reddit

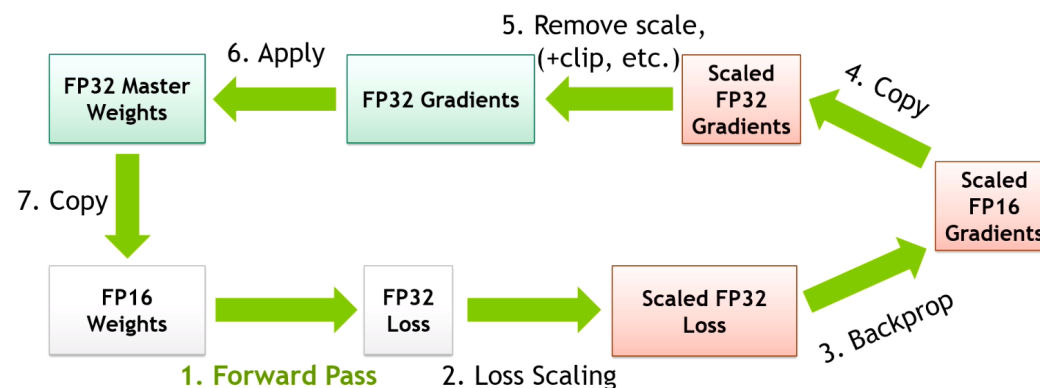


(d) multi-head SpMM on Yelp

# CogDL - Mixed Precision Training

- Support mixed precision training:
  - 1.44x~2.02x speedups due to half-precision (FP16) computation

## MIXED PRECISION TRAINING

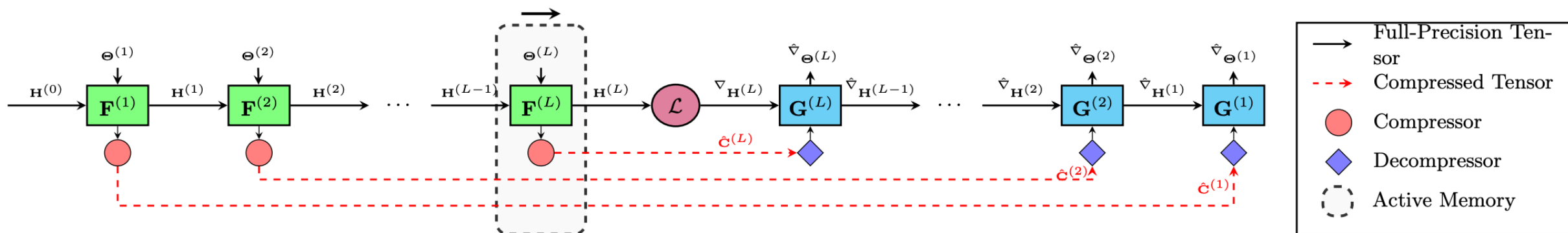


2-layer GCNs on the Reddit dataset ( single / mixed precision )

	GPU Memory	Accuracy	Training Speed	
			2080 Ti	3090
Single precision training	5,567 MB	95.44 ± 0.07	2.20 it/s	3.93 it/s
Mixed precison training	4,046 MB	95.35 ± 0.08	3.17 it/s	7.97 it/s

# CogDL - Activation Compressed Training

- Support activation compressed training:
  - reduce **6.4x~16x** training memory footprints

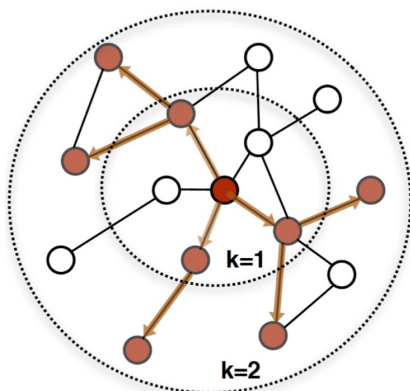


2-layer GCNs on three datasets: accuracy (%), act memory(MB)

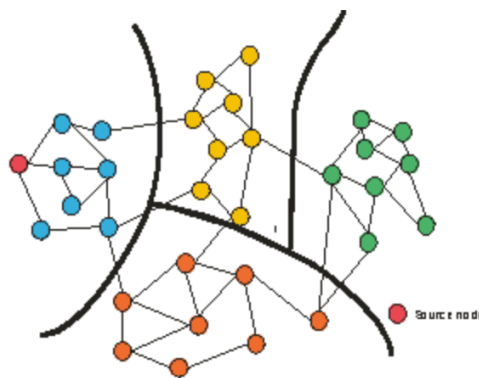
	origin (32-bit)	+actnn (4-bit)	+actnn (3-bit)	+actnn (2-bit)
Flickr	51.17 $\pm$ 0.19 (288)	51.08 $\pm$ 0.18 (37)	51.14 $\pm$ 0.18 (26)	51.20 $\pm$ 0.18 (18)
Reddit	95.33 $\pm$ 0.07 (1532)	95.32 $\pm$ 0.07 (194)	95.31 $\pm$ 0.07 (158)	95.34 $\pm$ 0.06 (112)
Yelp	39.86 $\pm$ 0.94 (4963)	40.06 $\pm$ 0.74 (773)	40.21 $\pm$ 0.82 (665)	39.89 $\pm$ 1.45 (551)

# CogDL – Training on Large-scale Graphs

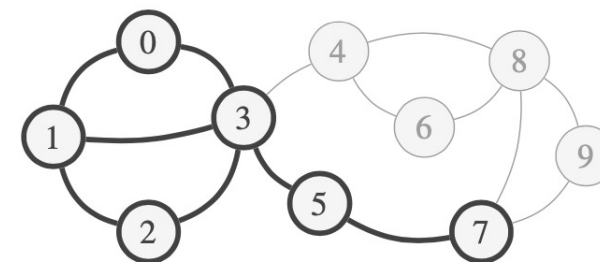
- Full-batch training on large-scale graphs is unaffordable
- Training GNNs via mini-batch sampling
- Usage: `python scripts/train.py --model gcn --dataset reddit --dw cluster_dw`



Neighbor Sampling  
(NeurIPS '17)



ClusterGCN  
(KDD '19)



$$\mathcal{G}_s = \text{SAMPLE}(\mathcal{G})$$

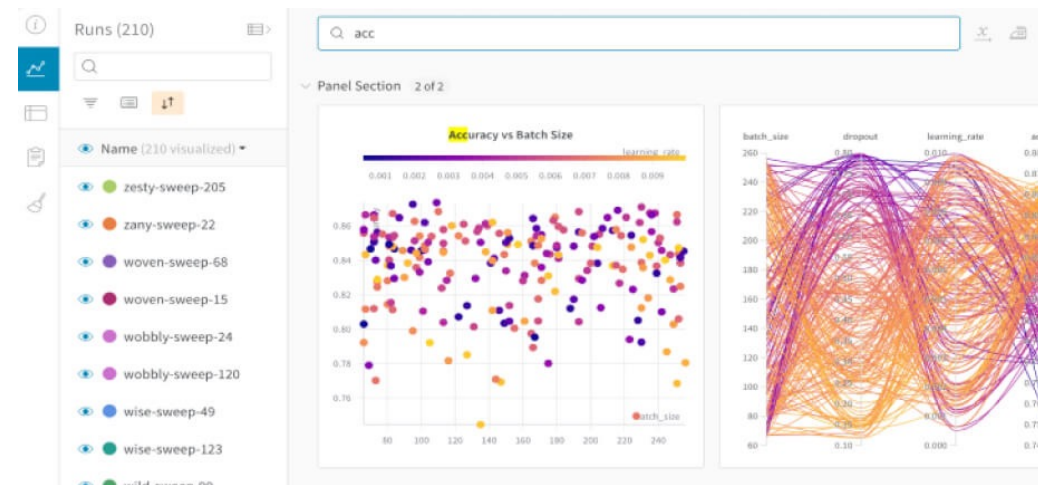
GraphSAINT  
(ICLR '20)

1. Will Hamilton, Zhitaoying, and Jure Leskovec. Inductive representation learning on large graphs. In NeurIPS '17.
2. Wei-Lin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, and Cho-Jui Hsieh. Cluster-GCN: An efficient algorithm for training deep and large graph convolutional networks. In KDD '19.
3. Hanqing Zeng, Hongkuan Zhou, Ajitesh Srivastava, Rajgopal Kannan, and Viktor Prasanna. Graphsaint: Graph sampling based inductive learning method. In ICLR '20.

# CogDL – Experiment Management

- Hyper-parameter Search
  - Integrate optuna for users to enable hyper-parameter search
- Experiment Management
  - Support Tensorboard and WandB for logging and debugging

```
def func(trial):  
    return {  
        "lr": trial.suggest_categorical("lr", [0.01, 0.1]),  
        "hidden_size": trial.suggest_categorical("hidden_size", [32, 64, 128]),  
        "dropout": trial.suggest_uniform("dropout", 0.1, 0.8),  
    }  
experiment(dataset="cora", model="gcn", search_space=func) # hyper-parameter search
```



# CogDL – Benchmarks

- Easy-to-reproduce benchmarks
  - Traditional graph tasks such as network embedding
  - Frontier graph benchmarks such as IGB<sup>[1]</sup>, GRB<sup>[2]</sup>, HGB<sup>[3]</sup>

Method	PPI (50%)	Wikipedia (50%)	Blogcatalog (50%)	DBLP (5%)	Flickr (5%)	<i>Reproducible</i>
NetMF [37]	23.73 $\pm$ 0.22	57.42 $\pm$ 0.56	42.47 $\pm$ 0.35	56.72 $\pm$ 0.14	36.27 $\pm$ 0.17	Yes
ProNE [66]	24.60 $\pm$ 0.39	56.06 $\pm$ 0.48	41.16 $\pm$ 0.26	56.85 $\pm$ 0.28	36.56 $\pm$ 0.11	Yes
NetSMF [36]	23.88 $\pm$ 0.35	53.81 $\pm$ 0.58	40.62 $\pm$ 0.35	59.76 $\pm$ 0.41	35.49 $\pm$ 0.07	Yes
Node2vec [17]	20.67 $\pm$ 0.54	54.59 $\pm$ 0.51	40.16 $\pm$ 0.29	57.36 $\pm$ 0.39	36.13 $\pm$ 0.13	Yes
LINE [45]	21.82 $\pm$ 0.56	52.46 $\pm$ 0.26	38.06 $\pm$ 0.39	49.78 $\pm$ 0.37	31.61 $\pm$ 0.09	Yes
DeepWalk [35]	20.74 $\pm$ 0.40	49.53 $\pm$ 0.54	40.48 $\pm$ 0.47	57.54 $\pm$ 0.32	36.09 $\pm$ 0.10	Yes
SpectralClustering [48]	22.48 $\pm$ 0.30	49.35 $\pm$ 0.34	41.41 $\pm$ 0.34	43.68 $\pm$ 0.58	33.09 $\pm$ 0.07	Yes
Hope [33]	21.43 $\pm$ 0.32	54.04 $\pm$ 0.47	33.99 $\pm$ 0.35	56.15 $\pm$ 0.22	28.97 $\pm$ 0.19	Yes
GraRep [5]	20.60 $\pm$ 0.34	54.37 $\pm$ 0.40	33.48 $\pm$ 0.30	52.76 $\pm$ 0.42	31.83 $\pm$ 0.12	Yes

1. Ziang Li, Ming Ding, Weikai Li, Zihan Wang, Ziyu Zeng, Yukuo Cen, and Jie Tang. Rethinking the Setting of Semi-supervised Learning on Graphs. In IJCAI'22.

2. Qinkai Zheng, Xu Zou, Yuxiao Dong, Yukuo Cen, Da Yin, Jiarong Xu, Yang Yang, and Jie Tang. Graph Robustness Benchmark: Benchmarking the Adversarial Robustness of Graph Machine Learning. In NeurIPS'21 D&B.

3. Qingsong Lv, Ming Ding, Qiang Liu, Yuxiang Chen, Wenzheng Feng, Siming He, Chang Zhou, Jian-guo Jiang, Yuxiao Dong, and Jie Tang. Are we really making much progress? Revisiting, benchmarking and refining the Heterogeneous Graph Neural Networks. In KDD'21.



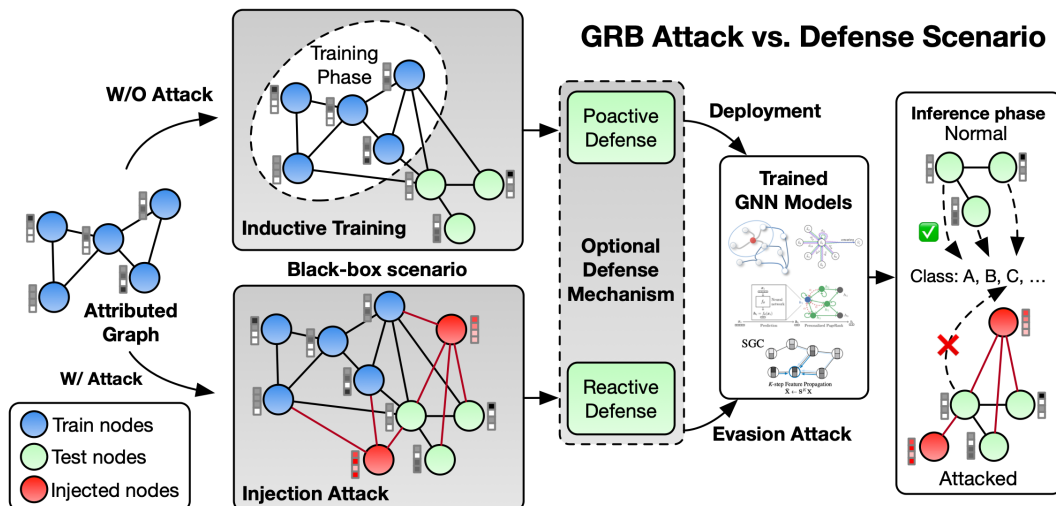
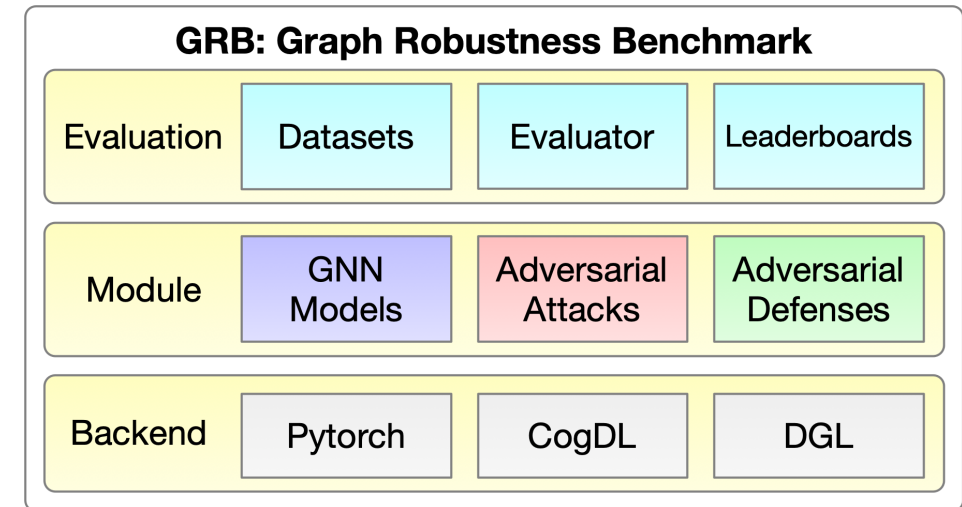
# CogDL – GRB

## Background:

Recently, works have proved that adversarial attacks can threaten the *robustness* of graph ML models in various tasks.

## Problems:

1. Ill-defined threat model in previous works.
2. Absence of unified and standard evaluation approach.



Example of GRB evaluation scenario

## Solution: Graph Robustness Benchmark (GRB)

*Scalable, general, unified, and reproducible* benchmark on *adversarial robustness* of graph ML models, which facilitates fair comparisons among various attacks & defenses and promotes future research in this field.

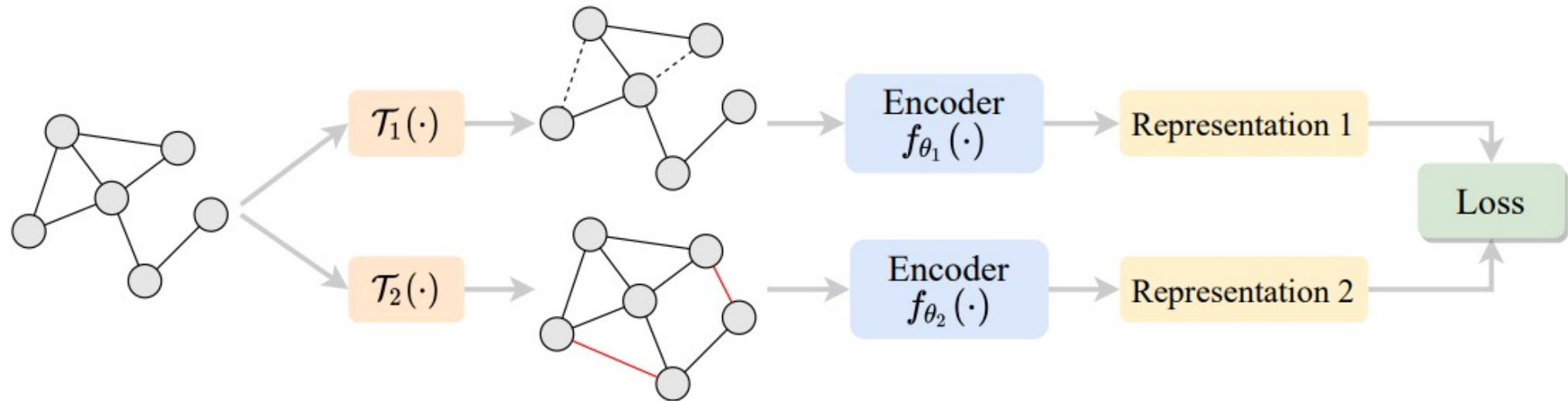
*Graph Robustness Benchmark: Rethinking and Benchmarking Adversarial Robustness of Graph Neural Networks*  
Qinkai Zheng, Xu Zou, Yuxiao Dong, Yukuo Cen, Jie Tang



# Contrastive Learning on Graphs

# Paradigm of Graph Contrastive Learning

- The graph contrastive learning
  - 1) generates two views based on different augmentations
  - 2) encodes the graphs of two views
  - 3) construct the self-supervised signal via contrast

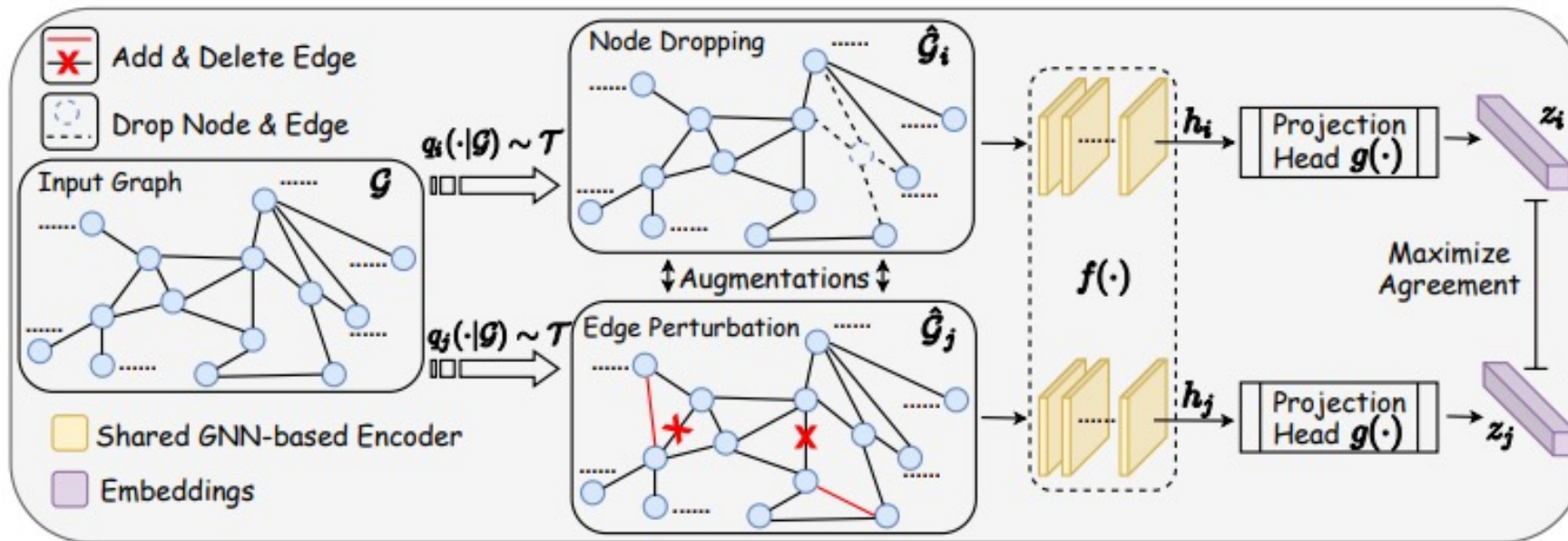


# Graph Contrastive Learning with Augmentations

GraphCL: a contrastive learning method with augmentations

- (1) propose different graph augmentation strategies
- (2) use projection heads for the graph encoding

- (3) maximizing the consistency 
$$\ell_n = -\log \frac{\exp(\text{sim}(\mathbf{z}_{n,i}, \mathbf{z}_{n,j})/\tau)}{\sum_{n'=1, n' \neq n}^N \exp(\text{sim}(\mathbf{z}_{n,i}, \mathbf{z}_{n',j})/\tau)}$$



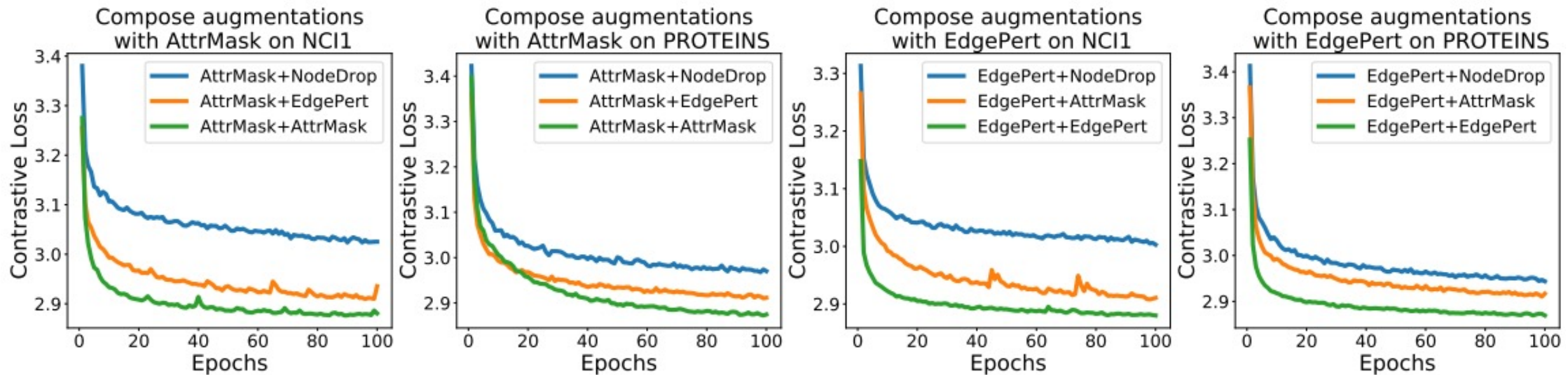
# Data Augmentation for Graphs

- Focus on three categories of graphs:
  - Biochemical molecules, social networks, image super-pixel graphs
- Propose four graph-level augmentations

Data augmentation	Type	Underlying Prior
Node dropping	Nodes, edges	Vertex missing does not alter semantics.
Edge perturbation	Edges	Semantic robustness against connectivity variations.
Attribute masking	Nodes	Semantic robustness against losing partial attributes.
Subgraph	Nodes, edges	Local structure can hint the full semantics.

# Role of Data Augmentation

- Contrastive loss curves for different augmentation pairs
- Using the same type gets better performance





# Experimental Results of GraphCL

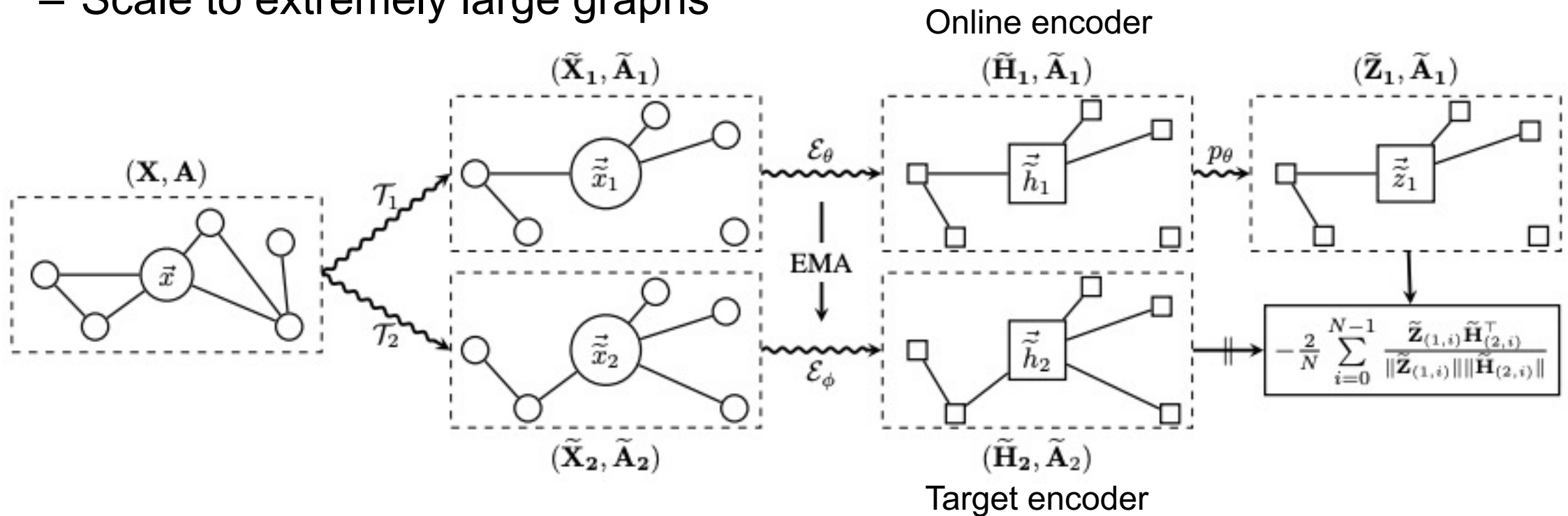
- Unsupervised representation learning setting

Datasets	Category	Graph Num.	Avg. Node	Avg. Degree
NCI1	Biochemical Molecules	4110	29.87	1.08
PROTEINS	Biochemical Molecules	1113	39.06	1.86
COLLAB	Social Networks	5000	74.49	32.99
RDT-B	Social Networks	2000	429.63	1.15

Dataset	NCI1	PROTEINS	DD	MUTAG	COLLAB	RDT-B	RDT-M5K	IMDB-B
GL	-	-	-	81.66±2.11	-	77.34±0.18	41.01±0.17	65.87±0.98
WL	80.01±0.50	72.92±0.56	-	80.72±3.00	-	68.82±0.41	46.06±0.21	72.30±3.44
DGK	80.31±0.46	73.30±0.82	-	87.44±2.72	-	78.04±0.39	41.27±0.18	66.96±0.56
node2vec	54.89±1.61	57.49±3.57	-	72.63±10.20	-	-	-	-
sub2vec	52.84±1.47	53.03±5.55	-	61.05±15.80	-	71.48±0.41	36.68±0.42	55.26±1.54
graph2vec	73.22±1.81	73.30±2.05	-	83.15±9.25	-	75.78±1.03	47.86±0.26	71.10±0.54
InfoGraph	76.20±1.06	74.44±0.31	72.85±1.78	89.01±1.13	70.65±1.13	82.50±1.42	53.46±1.03	73.03±0.87
GraphCL	77.87±0.41	74.39±0.45	78.62±0.40	86.80±1.34	71.36±1.15	89.53±0.84	55.99±0.28	71.14±0.44

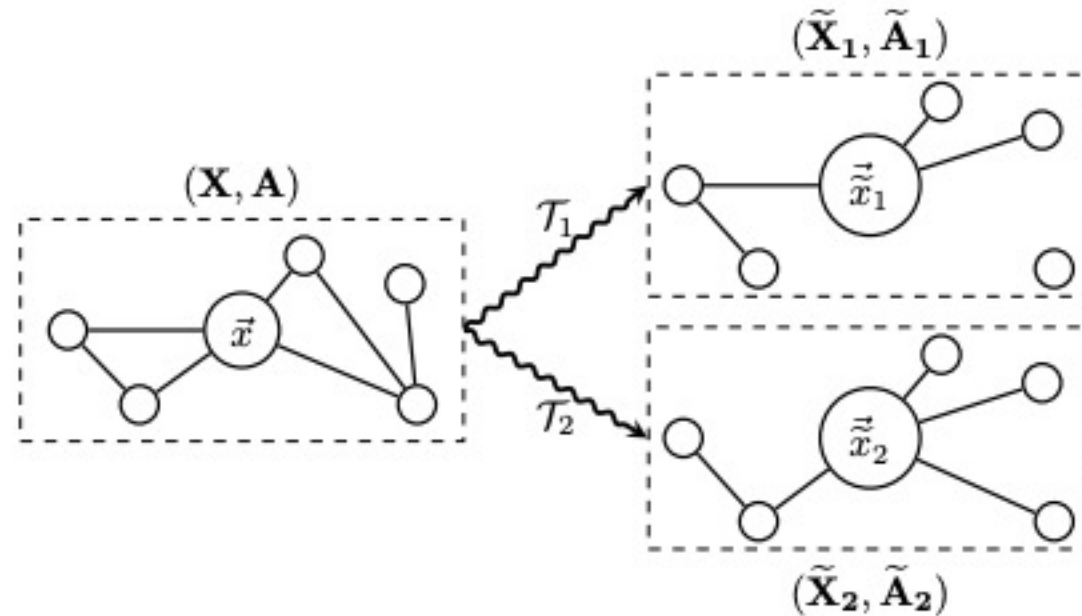
# Bootstrapped Graph Latents (BGRL)

- The characteristics of BGRL:
  - Simple graph augmentation
  - Not requiring negative samples
  - Scale to extremely large graphs



# Graph Augmentation of BGRL

- Applying stochastic graph augmentation functions
  - random node feature masking: Bernoulli distribution  $\mathcal{B}(1 - p_f)$
  - random edge masking: Bernoulli distribution  $\mathcal{B}(1 - p_e)$
- With fixed hyperparameters for each graph





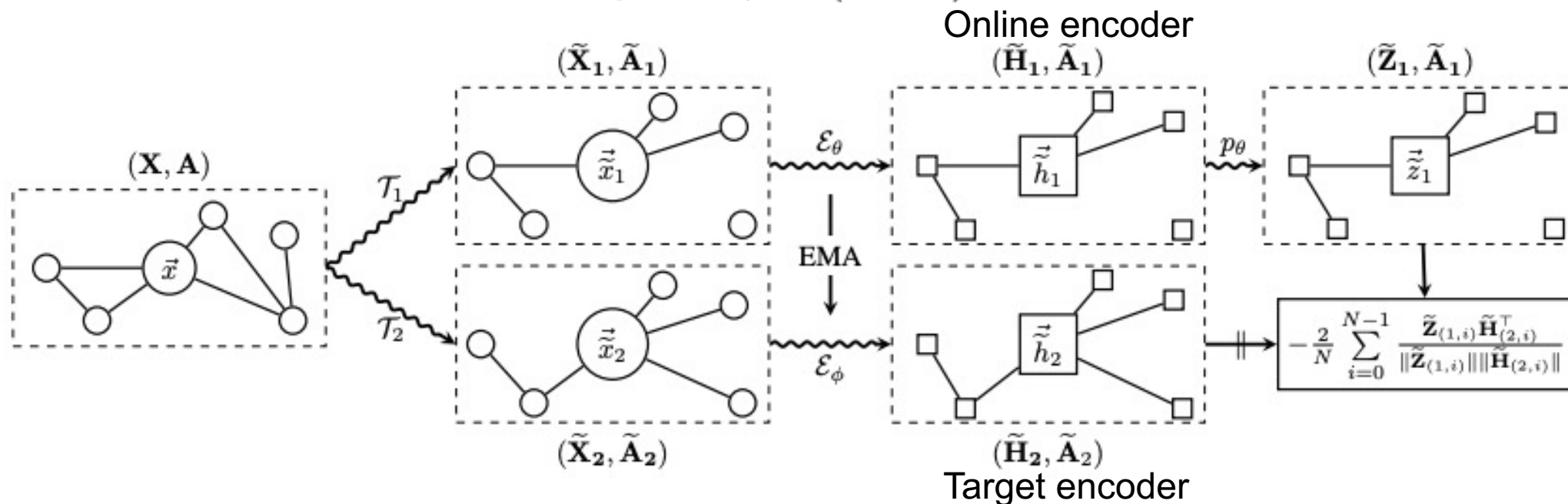
# BGRL Update Step

- Update the online encoder via gradient descent

$$\ell(\theta, \phi) = -\frac{2}{N} \sum_{i=0}^{N-1} \frac{\tilde{\mathbf{Z}}_{(1,i)} \tilde{\mathbf{H}}_{(2,i)}^\top}{\|\tilde{\mathbf{Z}}_{(1,i)}\| \|\tilde{\mathbf{H}}_{(2,i)}\|}$$

- Update the target encoder via EMA

$$\phi \leftarrow \tau\phi + (1 - \tau)\theta$$



# Experimental Results of BGRL

	Task	Nodes	Edges	Features	Classes
<b>WikiCS</b>	Transductive	11,701	216,123	300	10
<b>Amazon Computers</b>	Transductive	13,752	245,861	767	10
<b>Amazon Photos</b>	Transductive	7,650	119,081	745	8
<b>Coauthor CS</b>	Transductive	18,333	81,894	6,805	15
<b>Coauthor Physics</b>	Transductive	34,493	247,962	8,415	5

	WikiCS	Am. Comp.	Am. Photos	Co.CS	Co.Phy
Raw features	71.98 $\pm$ 0.00	73.81 $\pm$ 0.00	78.53 $\pm$ 0.00	90.37 $\pm$ 0.00	93.58 $\pm$ 0.00
DeepWalk	74.35 $\pm$ 0.06	85.68 $\pm$ 0.06	89.44 $\pm$ 0.11	84.61 $\pm$ 0.22	91.77 $\pm$ 0.15
DeepWalk + feat.	77.21 $\pm$ 0.03	86.28 $\pm$ 0.07	90.05 $\pm$ 0.08	87.70 $\pm$ 0.04	94.90 $\pm$ 0.09
DGI	75.35 $\pm$ 0.14	83.95 $\pm$ 0.47	91.61 $\pm$ 0.22	92.15 $\pm$ 0.63	94.51 $\pm$ 0.52
GMI	74.85 $\pm$ 0.08	82.21 $\pm$ 0.31	90.68 $\pm$ 0.17	OOM	OOM
MVGRL	77.52 $\pm$ 0.08	87.52 $\pm$ 0.11	91.74 $\pm$ 0.07	92.11 $\pm$ 0.12	95.33 $\pm$ 0.03
Random-Init*	78.95 $\pm$ 0.58	86.46 $\pm$ 0.38	92.08 $\pm$ 0.48	91.64 $\pm$ 0.29	93.71 $\pm$ 0.29
GRACE *	<b>80.14 <math>\pm</math> 0.48</b>	89.53 $\pm$ 0.35	92.78 $\pm$ 0.45	91.12 $\pm$ 0.20	OOM
BGRL*	79.98 $\pm$ 0.10	<b>90.34 <math>\pm</math> 0.19</b>	<b>93.17 <math>\pm</math> 0.30</b>	<b>93.31 <math>\pm</math> 0.13</b>	<b>95.73 <math>\pm</math> 0.05</b>
GCA	78.35 $\pm$ 0.05	88.94 $\pm$ 0.15	92.53 $\pm$ 0.16	93.10 $\pm$ 0.01	95.73 $\pm$ 0.03
Supervised GCN	77.19 $\pm$ 0.12	86.51 $\pm$ 0.54	92.42 $\pm$ 0.22	93.03 $\pm$ 0.31	95.65 $\pm$ 0.16

# Experimental Results of BGRL (cont.)

- Comparison with SOTA methods on ogbn-arxiv and PPI datasets

	Validation	Test
MLP	$57.65 \pm 0.12$	$55.50 \pm 0.23$
node2vec	$71.29 \pm 0.13$	$70.07 \pm 0.13$
Random-Init*	$69.90 \pm 0.11$	$68.94 \pm 0.15$
DGI*	$71.26 \pm 0.11$	$70.34 \pm 0.16$
GRACE full-graph*	OOM	OOM
GRACE-SUBSAMPLING ( $k = 2$ )*	$60.49 \pm 3.72$	$60.24 \pm 4.06$
GRACE-SUBSAMPLING ( $k = 8$ )*	$71.30 \pm 0.17$	$70.33 \pm 0.18$
GRACE-SUBSAMPLING ( $k = 32$ )*	$72.18 \pm 0.16$	$71.18 \pm 0.16$
GRACE-SUBSAMPLING ( $k = 2048$ )*	<b><math>72.61 \pm 0.15</math></b>	<b><math>71.51 \pm 0.11</math></b>
BGRL*	<b><math>72.53 \pm 0.09</math></b>	<b><math>71.64 \pm 0.12</math></b>
Supervised GCN	$73.00 \pm 0.17$	$71.74 \pm 0.29$

ogbn-arxiv

	PPI
Raw features	42.20
DGI	$63.80 \pm 0.20$
GMI	$65.00 \pm 0.02$
Random-Init	$62.60 \pm 0.20$
GRACE MeanPooling encoder*	$69.66 \pm 0.15$
BGRL MeanPooling encoder*	$69.41 \pm 0.15$
GRACE GAT encoder*	$69.71 \pm 0.17$
BGRL GAT encoder*	<b><math>70.49 \pm 0.05</math></b>
Supervised MeanPooling	$96.90 \pm 0.20$
Supervised GAT	$97.30 \pm 0.20$

PPI

# Experiment on MAG240M

- MAG240M: over 240 million nodes (of which 1.4 million are labeled) and 1.7 billion edges

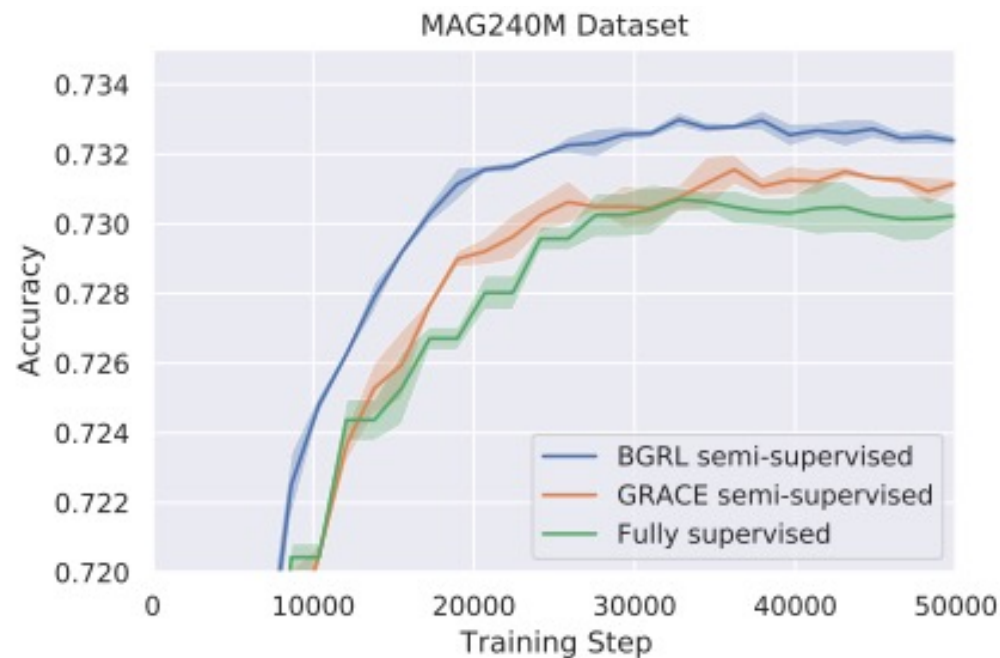


Figure 4: Performance on MAG240M using BGRL or GRACE-SUBSAMPLING as an auxiliary signal, averaged over 5 seeds and run for 50k steps.

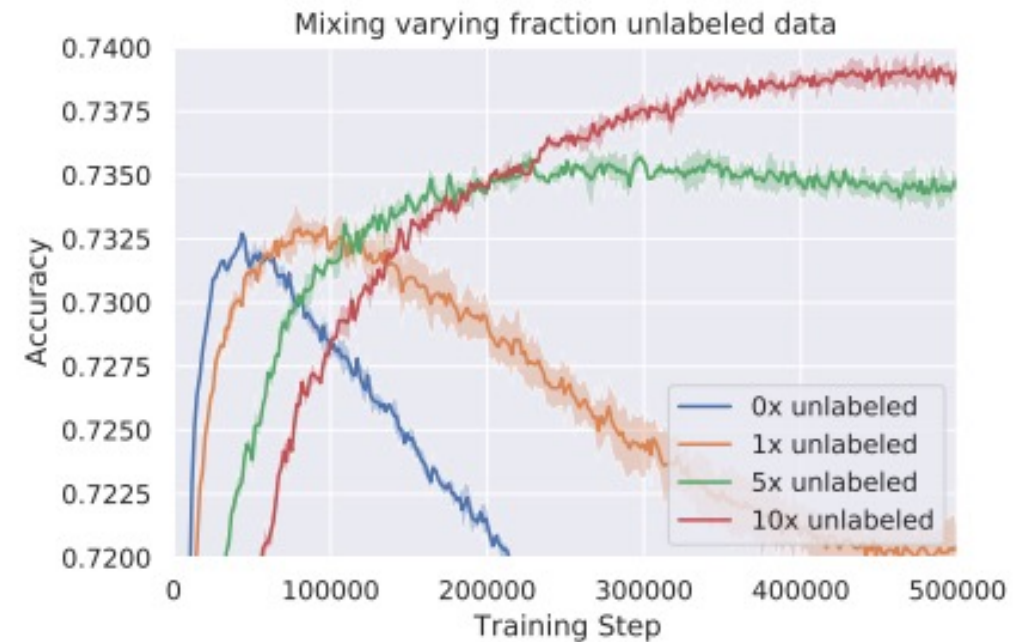
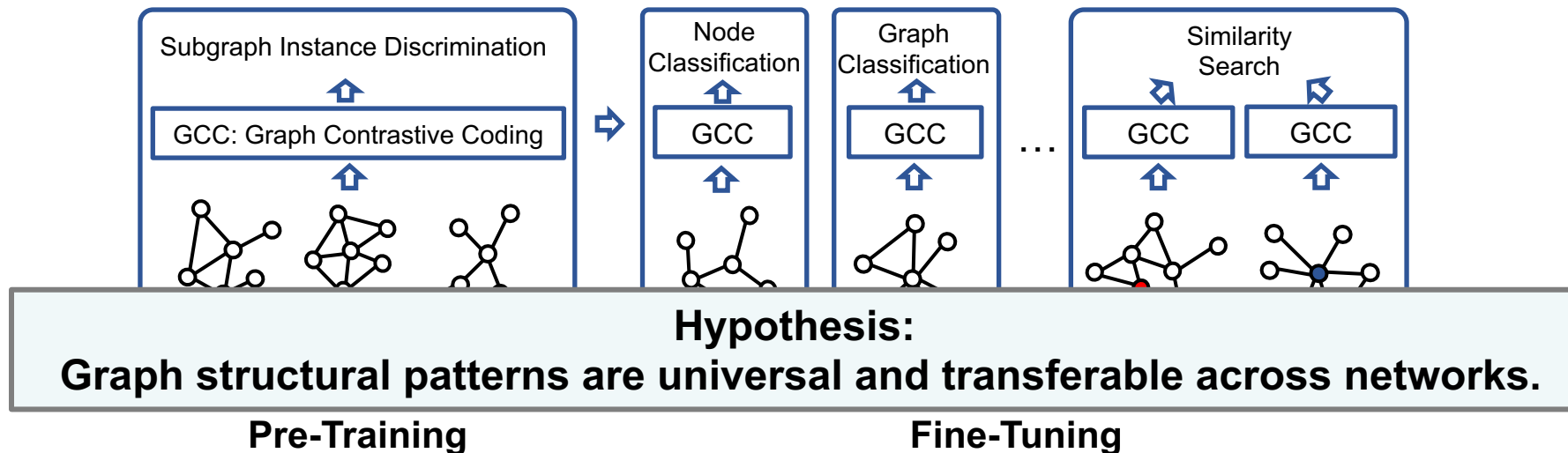


Figure 5: Mixing varying amounts of unlabeled data for representation learning with BGRL, averaged over 5 seeds and run for 500k steps.

# Graph Contrastive Coding (GCC)

- Problem:
  - Learn a function  $f$  that maps a vertex to a low-dimensional vector
  - **Structural similarity**: map vertices with similar local network topologies close in the vector space
  - **Transferability**: compatible with vertices and graphs from various sources, even unseen during training time.



# GCC Pre-training

- **Pre-training Task:** **Instance** Discrimination
- **InfoNCE Loss:** output **instance representations** that are capable of capturing the **similarities** between instances

$$\mathcal{L} = -\log \frac{\exp(\mathbf{q}^\top \mathbf{k}_+ / \tau)}{\sum_{i=0}^K \exp(\mathbf{q}^\top \mathbf{k}_i / \tau)}$$

- Contrastive learning for graphs?
  - **Q1:** How to define **instances** in graphs?
  - **Q2:** How to define **(dis) similar instance** pairs?
  - **Q3:** What are the proper **encoders**?

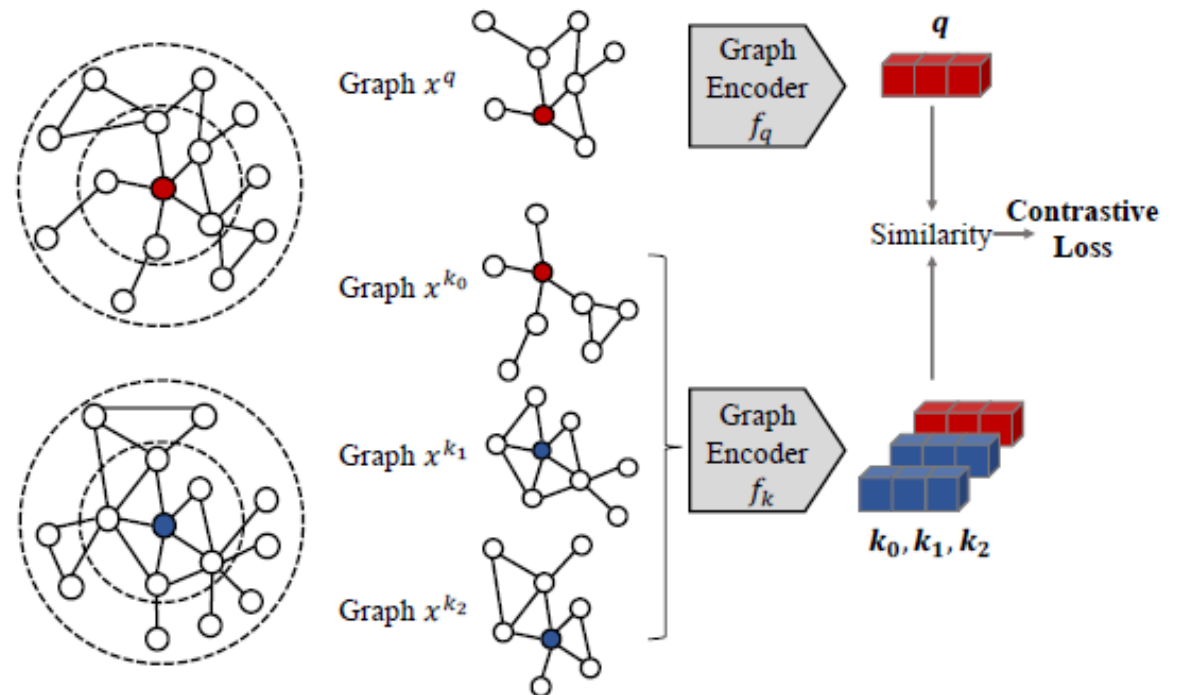
- query instance  $x^q$
- query  $\mathbf{q}$  (embedding of  $x^q$ ), i.e.,  $\mathbf{q} = f(x^q)$
- dictionary of keys  $\{\mathbf{k}_0, \mathbf{k}_1, \dots, \mathbf{k}_K\}$
- key  $\mathbf{k} = f(x^k)$



# GCC Pre-training

- **Q1**: How to define **instances** in graphs?
- **Q2**: How to define **(dis) similar instance**?
- **Q3**: What are the proper **encoders**?

$$\mathcal{L} = -\log \frac{\exp(\mathbf{q}^\top \mathbf{k}_+ / \tau)}{\sum_{i=0}^K \exp(\mathbf{q}^\top \mathbf{k}_i / \tau)}$$



# GCC Pre-training: Learning Algorithms

- Optimizing Contrastive Loss
  - Encoded query  $q$
  - $K + 1$  encoded keys  $\{k_0, \dots, k_K\}$

$$\mathcal{L} = -\log \frac{\exp(\mathbf{q}^\top \mathbf{k}_+ / \tau)}{\sum_{i=0}^K \exp(\mathbf{q}^\top \mathbf{k}_i / \tau)}$$

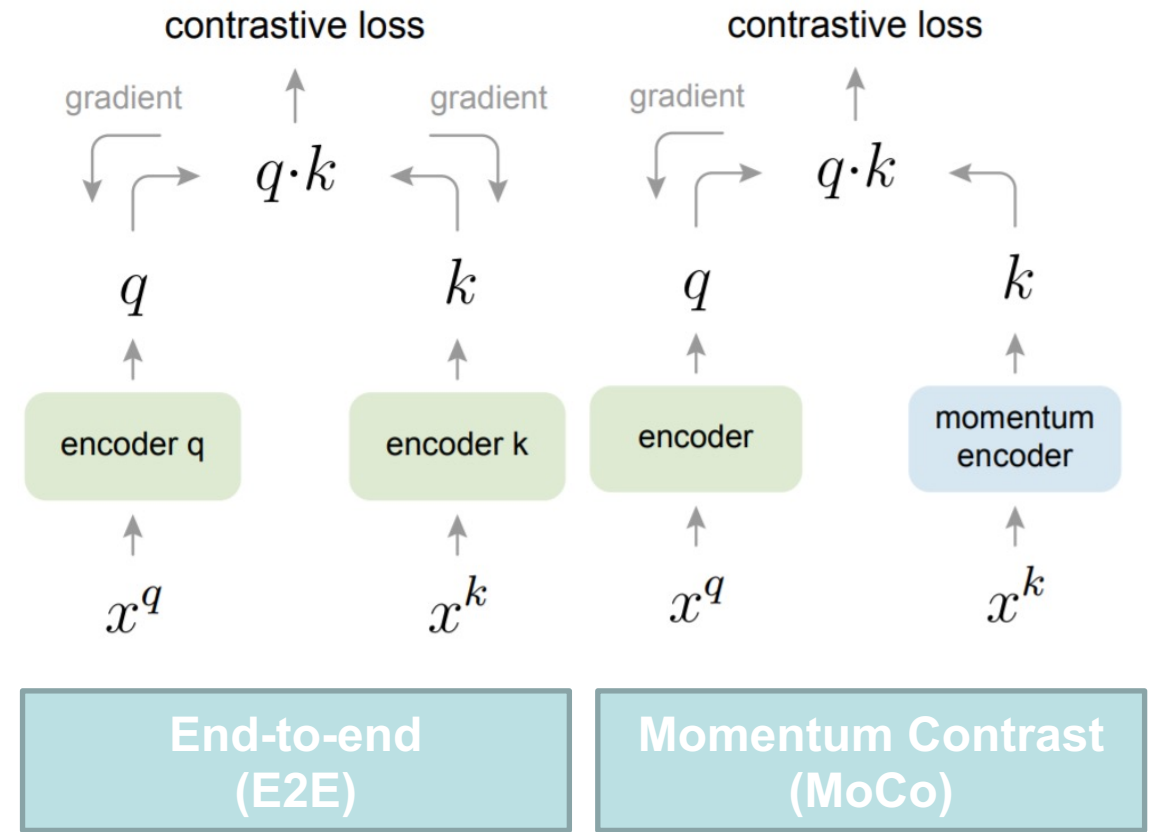
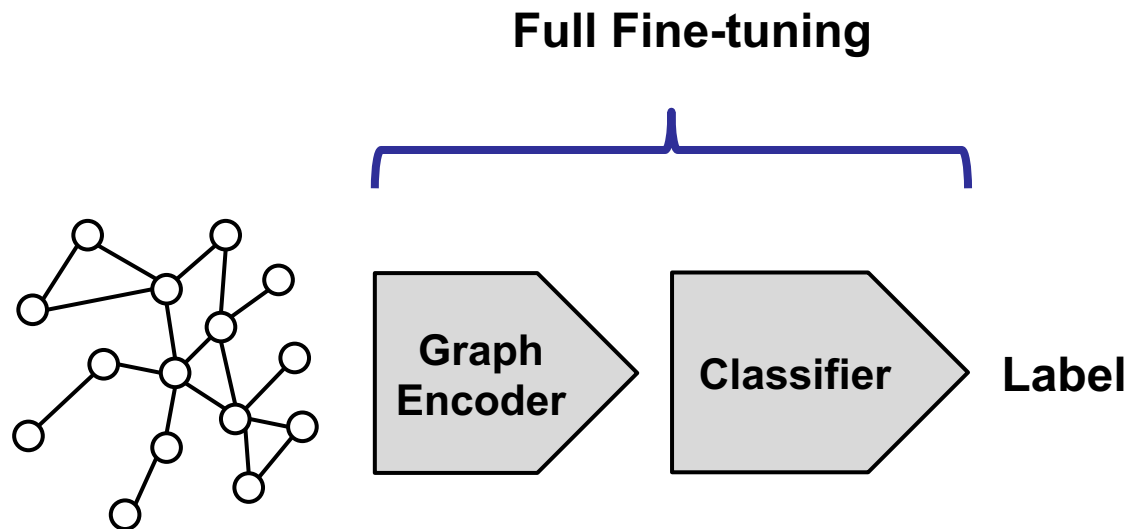


Figure Credit: Momentum Contrast for Unsupervised Visual Representation Learning ([arxiv.org/abs/1911.05722](https://arxiv.org/abs/1911.05722))

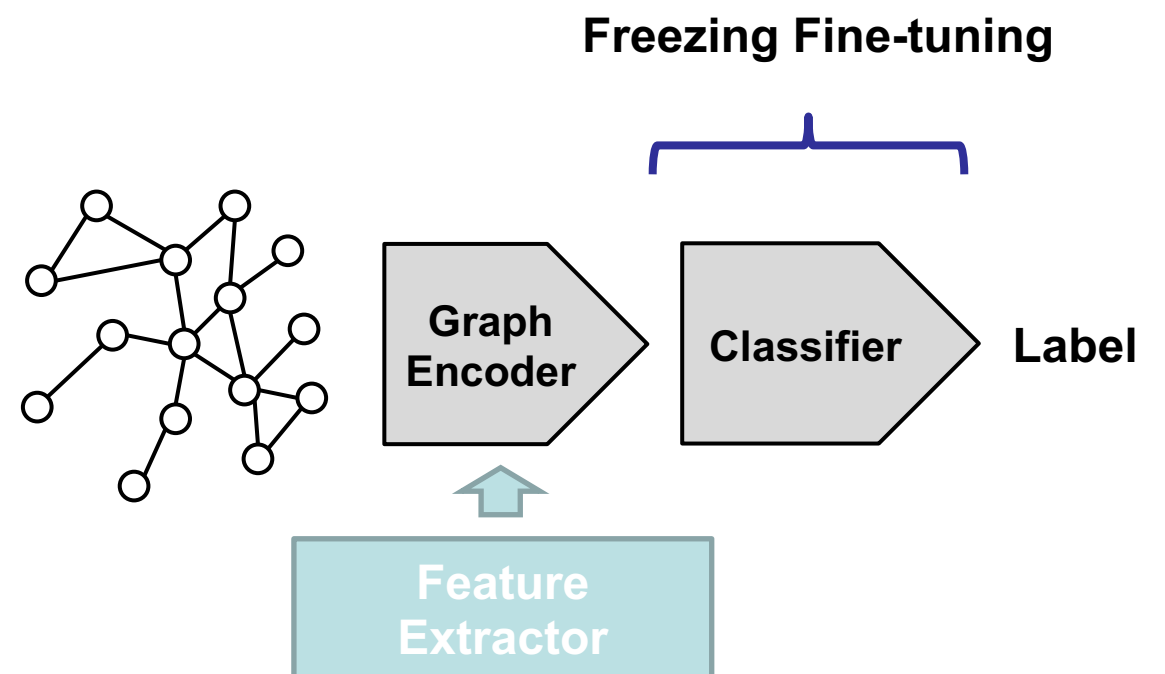


# GCC Fine-tuning: Full v.s. Freezing

## Full fine-tuning



## Freezing fine-tuning



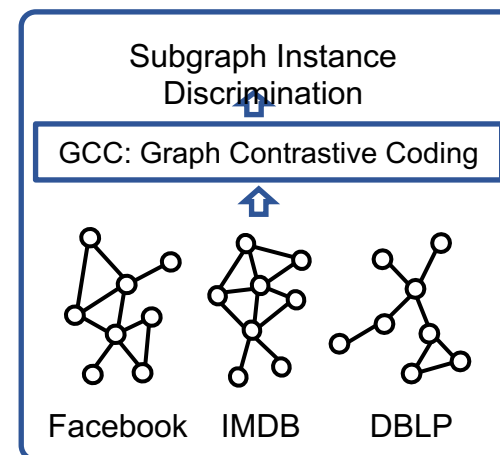
# GCC Pre-Training / Fine-tuning

- Six real-world information networks for pre-training.

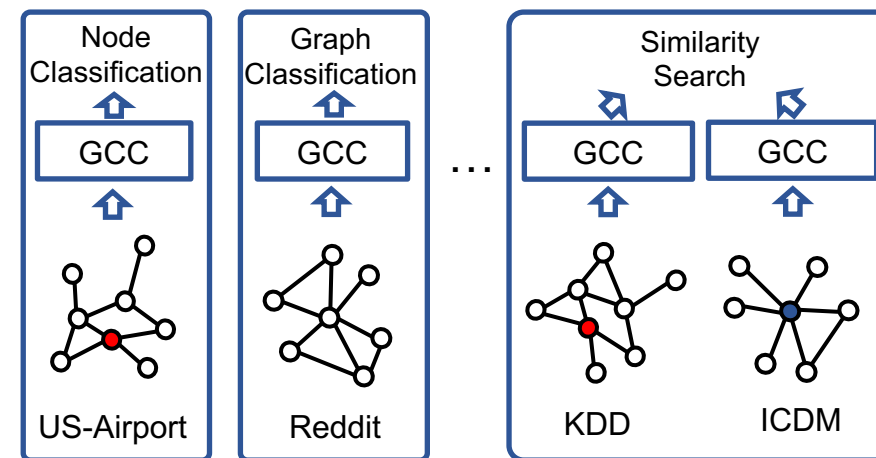
Table 1: Datasets for pre-training, sorted by number of vertices.

Dataset	Academia	DBLP (SNAP)	DBLP (NetRep)	IMDB	Facebook	LiveJournal
$ V $	137,969	317,080	540,486	896,305	3,097,165	4,843,953
$ E $	739,384	2,099,732	30,491,458	7,564,894	47,334,788	85,691,368

- Fine-tuning Tasks:
  - Node classification
  - Graph classification
  - Top-k Similarity search



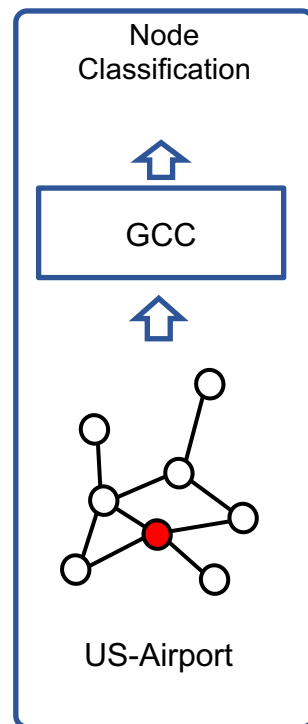
**Pre-Training**



**Fine-Tuning**

# Task 1: Node Classification

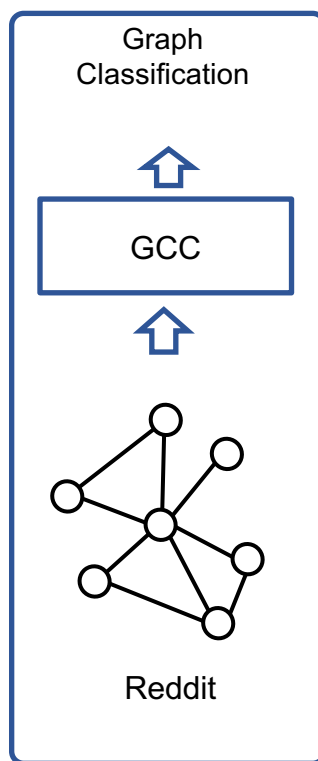
- Setup
  - US-Airport
  - AMiner academic graph



Datasets	US-Airport	H-index
$ V $	1,190	5,000
$ E $	13,599	44,020
ProNE	62.3	69.1
GraphWave	60.2	70.3
Struc2vec	<b>66.2</b>	> 1 Day
GCC (E2E, freeze)	64.8	<b>78.3</b>
GCC (MoCo, freeze)	65.6	75.2
GCC (rand, full)	64.2	76.9
GCC (E2E, full)	<b>68.3</b>	80.5
GCC (MoCo, full)	67.2	<b>80.6</b>

# Task 2: Graph Classification

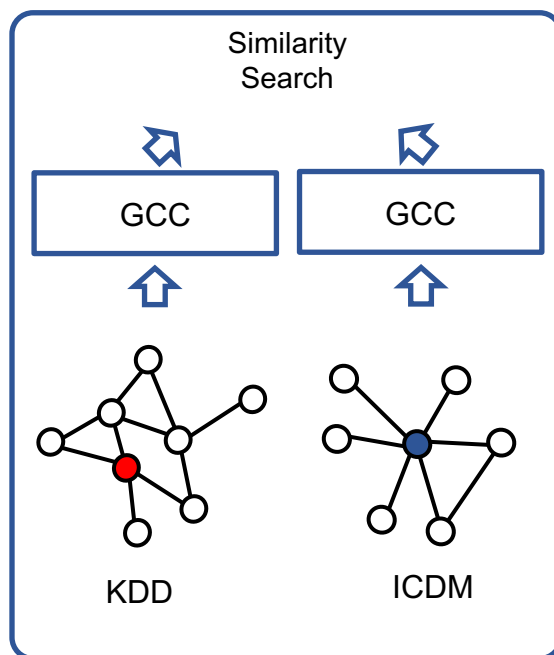
- Setup
  - COLLAB, RDT-B, RDT-M, & IMDB-B, IMDB-M



Datasets	IMDB-B	IMDB-M	COLLAB	RDT-B	RDT-M
# graphs	1,000	1,500	5,000	2,000	5,000
# classes	2	3	3	2	5
Avg. # nodes	19.8	13.0	74.5	429.6	508.5
DGK	67.0	44.6	73.1	78.0	41.3
graph2vec	71.1	50.4	–	75.8	47.9
InfoGraph	<b>73.0</b>	<b>49.7</b>	–	82.5	53.5
GCC (E2E, freeze)	71.7	49.3	74.7	87.5	52.6
GCC (MoCo, freeze)	72.0	49.4	<b>78.9</b>	<b>89.8</b>	<b>53.7</b>
DGCNN	70.0	47.8	73.7	–	–
GIN	<b>75.6</b>	<b>51.5</b>	80.2	<b>89.4</b>	<b>54.5</b>
GCC (rand, full)	<b>75.6</b>	50.9	79.4	87.8	52.1
GCC (E2E, full)	70.8	48.5	79.0	86.4	47.4
GCC (MoCo, full)	73.8	50.3	<b>81.1</b>	87.6	53.0

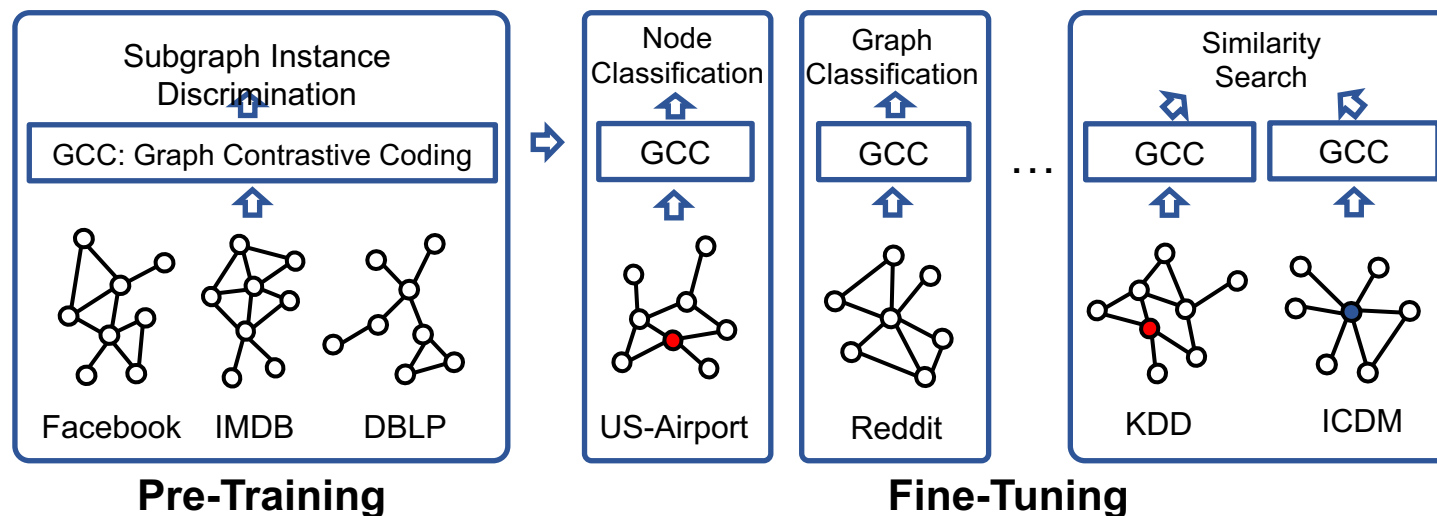
# Task 3: Top-k Similarity Search

- Setup
  - AMiner academic graph



	KDD-ICDM		SIGIR-CIKM		SIGMOD-ICDE	
$ V $	2,867	2,607	2,851	3,548	2,616	2,559
$ E $	7,637	4,774	6,354	7,076	8,304	6,668
# ground truth		697		874		898
$k$	20	40	20	40	20	40
Random	0.0198	0.0566	0.0223	0.0447	0.0221	0.0521
RoIX	0.0779	0.1288	0.0548	0.0984	0.0776	0.1309
Panther++	0.0892	0.1558	<b>0.0782</b>	0.1185	0.0921	0.1320
GraphWave	0.0846	<b>0.1693</b>	0.0549	0.0995	<b>0.0947</b>	<b>0.1470</b>
GCC (E2E)	<b>0.1047</b>	0.1564	0.0549	<b>0.1247</b>	0.0835	0.1336
GCC (MoCo)	0.0904	0.1521	0.0652	0.1178	0.0846	0.1425

# Summary of GCC



- We study the pre-training of GNN with the goal of characterizing and transferring structural representations in social and information networks.
- We present Graph Contrastive Coding, which is a graph-based contrastive learning framework to pre-train GNN.
- The pre-trained GNN achieves competitive performance to its supervised trained-from-scratch counterparts in 3 graph learning tasks on 10 graph datasets.

# Generative Learning on Graphs



# Web-Scale Graphs



Academic Graph



Microsoft Office Graph



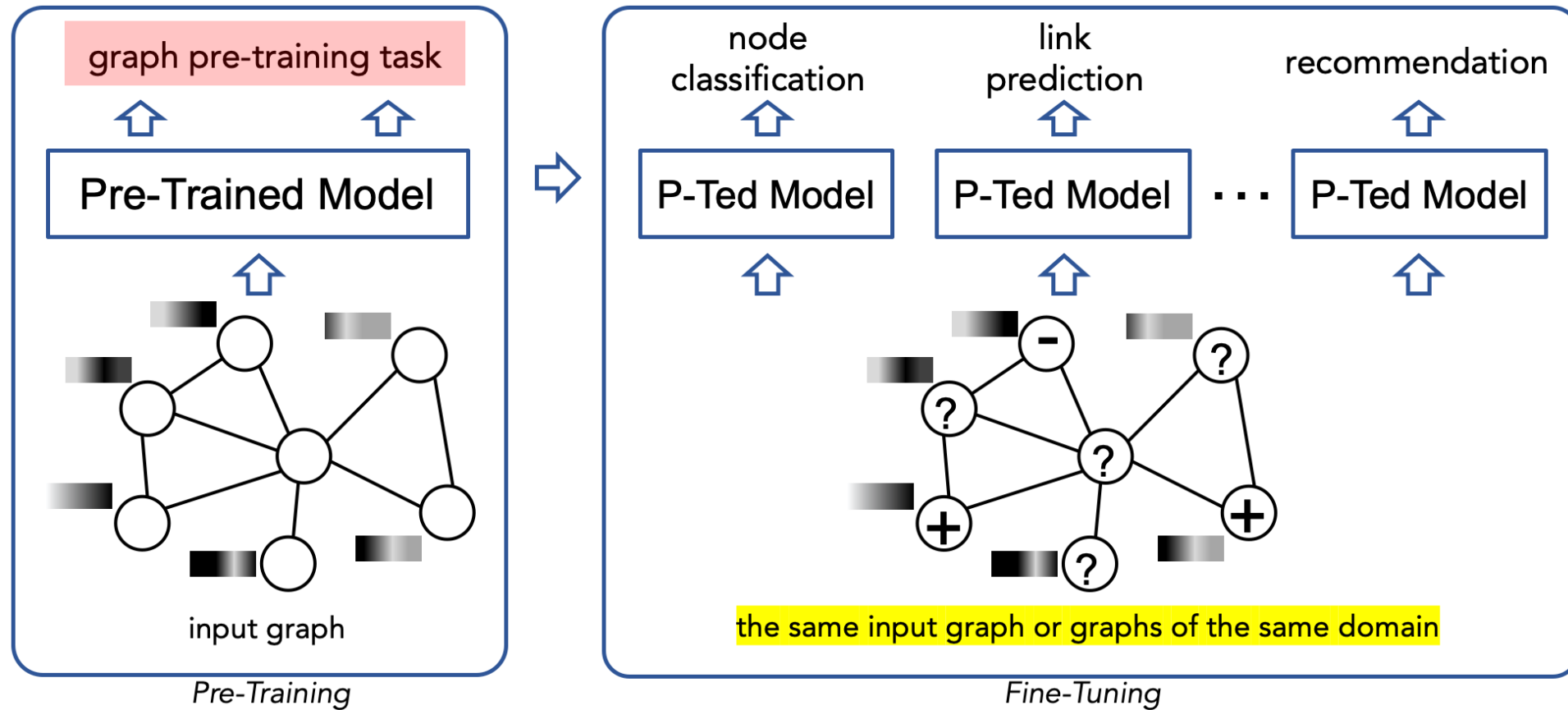
LinkedIn Economic Graph



Facebook Entity Graph



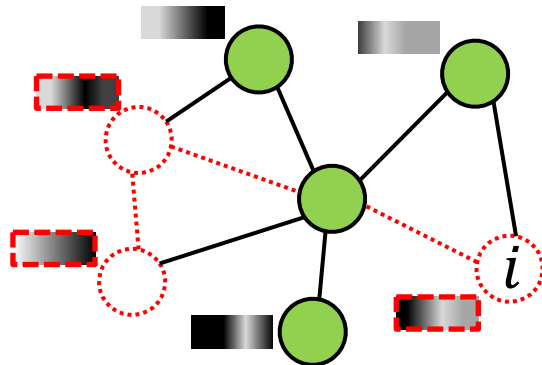
# GNN Pre-Training



# GPT-GNN: Generative Pre-Training of GNNs

- Model the graph distribution  $p(G; \theta)$  by learning to reconstruct the input graph.
  - Factorize the graph likelihood into two terms:
    - Attribute Generation
    - Edge Generation

$$\log p_{\theta}(X, E) = \sum_{i=1}^{|\mathcal{V}|} \log p_{\theta}(X_i, E_i | X_{<i}, E_{<i}).$$



attribute and edge **masked**  
input graph

$$p_{\theta}(X_i, E_i | X_{<i}, E_{<i}) \\ = p_{\theta}(X_i | X_{<i}, E_{<i}) \cdot p_{\theta}(E_i | X_{<i}, E_{<i})$$

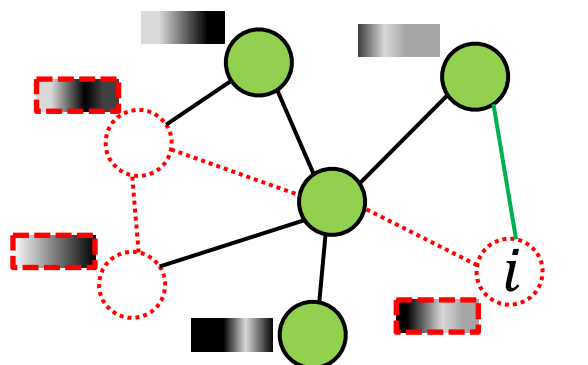
?

Lose the dependency between  $X_i$  and  $E_i$

# GPT-GNN: Generative Pre-Training of GNNs

- Model the graph distribution  $p(G; \theta)$  by learning to reconstruct the input graph.
  - Factorize the graph likelihood into two terms:
    - Attribute Generation: given observed edges, generate node attributes
    - Edge Generation: given observed edges and generated attributes, generate masked edges

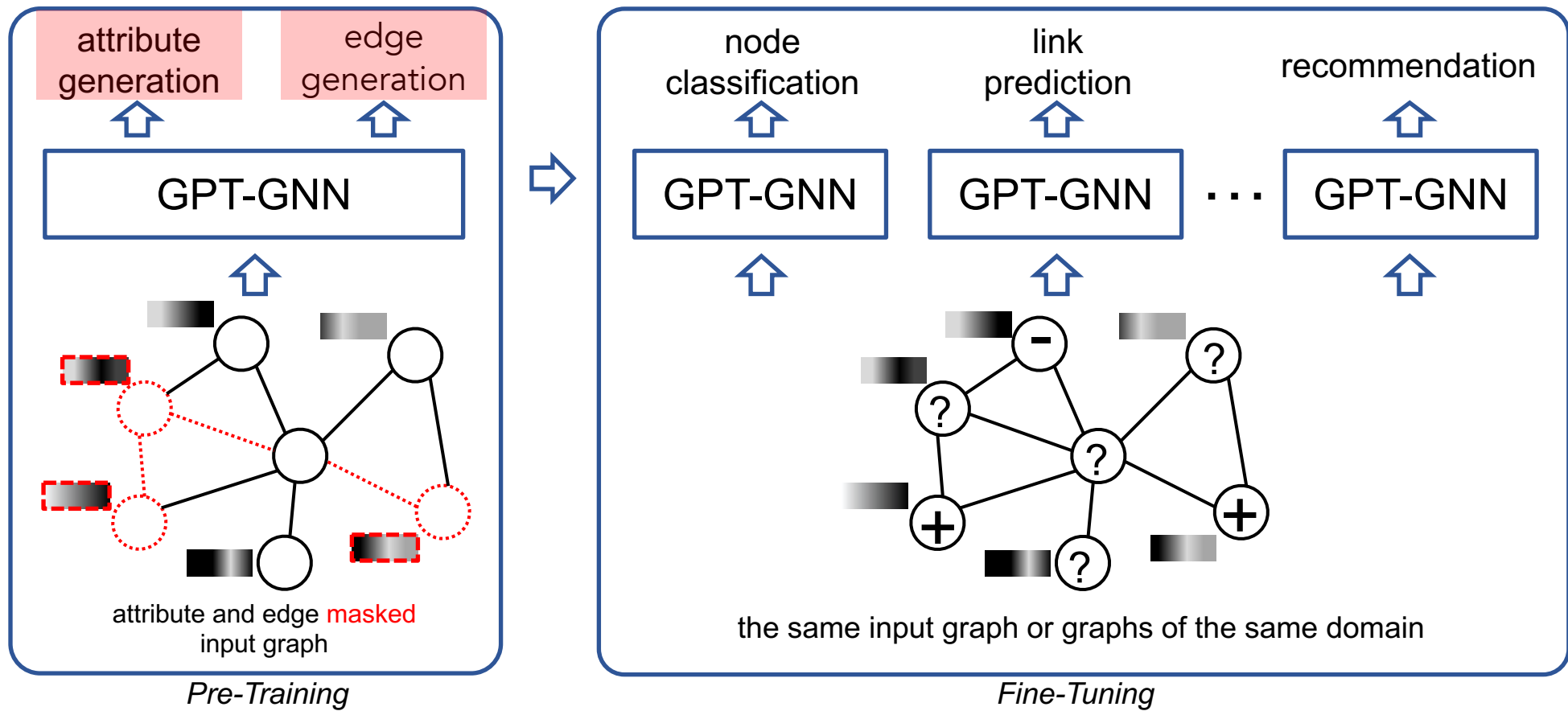
$$\log p_{\theta}(X, E) = \sum_{i=1}^{|\mathcal{V}|} \log p_{\theta}(X_i, E_i \mid X_{<i}, E_{<i}).$$



attribute and edge **masked**  
input graph

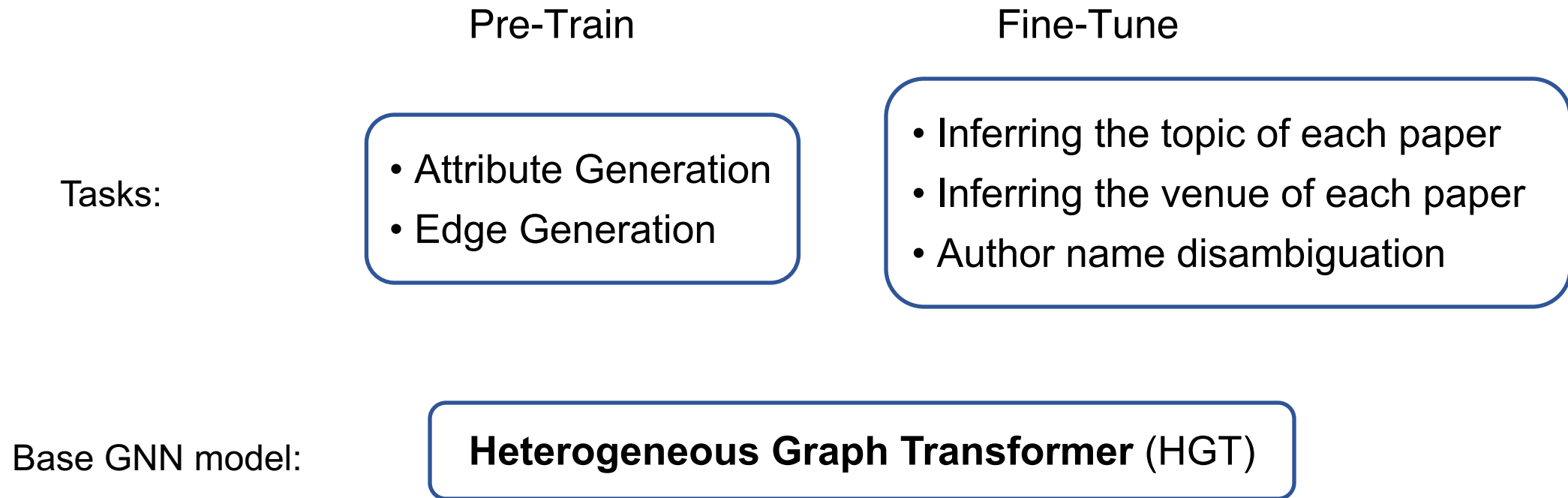
$$\begin{aligned} & p_{\theta}(X_i, E_i \mid X_{<i}, E_{<i}) \\ &= \sum_o p_{\theta}(X_i, E_{i,\neg o} \mid E_{i,o}, X_{<i}, E_{<i}) \cdot p_{\theta}(E_{i,o} \mid X_{<i}, E_{<i}) \\ &= \mathbb{E}_o \left[ p_{\theta}(X_i, E_{i,\neg o} \mid E_{i,o}, X_{<i}, E_{<i}) \right] \\ &= \mathbb{E}_o \left[ \underbrace{p_{\theta}(X_i \mid E_{i,o}, X_{<i}, E_{<i})}_{1) \text{ generate attributes}} \cdot \underbrace{p_{\theta}(E_{i,\neg o} \mid E_{i,o}, X_{\leq i}, E_{<i})}_{2) \text{ generate edges}} \right]. \end{aligned}$$

# GPT-GNN: Generative Pre-Training of GNNs



# GPT-GNN: Generative Pre-Training of GNNs

- Data: Microsoft Academic Graph



# GPT-GNN: Generative Pre-Training of GNNs

- Data: Microsoft Academic Graph

	Pre-Train	Fine-Tune
No Transfer:	CS Academic Graph	CS Academic Graph
Field Transfer:	Med, Bio, Physics...	CS Academic Graph
Time Transfer:	CS before 2014	CS after 2014
Time + Field Transfer:	Med, Bio, Physics... before 2014	CS after 2014

# GPT-GNN: Generative Pre-Training of GNNs

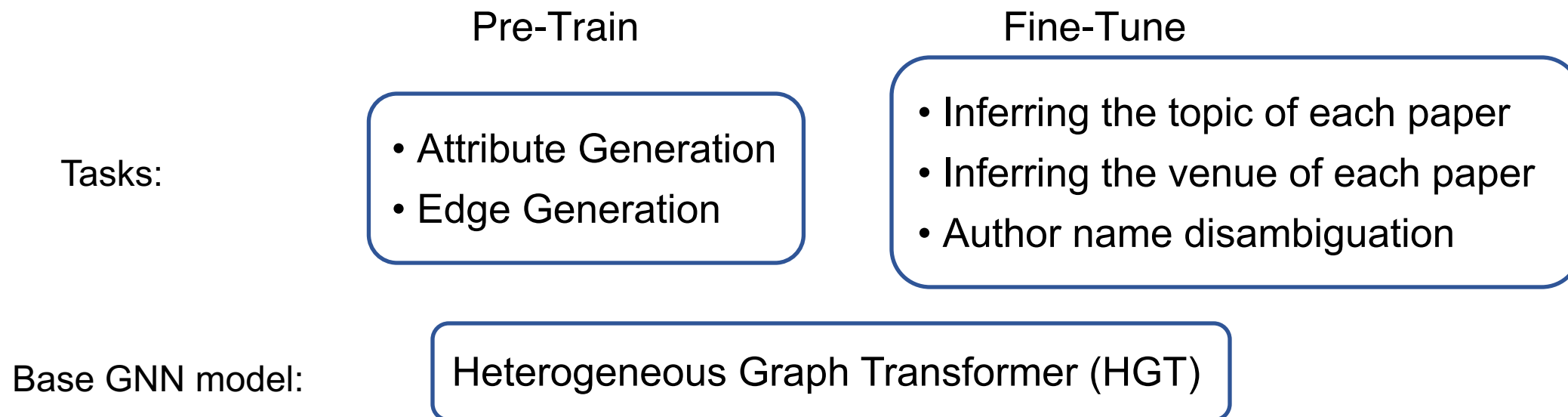
Downstream Dataset		OAG		
Evaluation Task		Paper-Field	Paper-Venue	Author ND
No Pre-train		.346±.149	.598±.122	.813±.105
Field Transfer	GAE	.403±.114	.626±.093	.836±.084
	GraphSAGE (unsp.)	.368±.125	.609±.096	.818±.092
	Graph Infomax	.387±.112	.612±.097	.827±.084
	GPT-GNN (Attr)	.396±.118	.623±.105	.834±.086
	GPT-GNN (Edge)	.413±.109	.635±.096	.842±.093
	GPT-GNN	<b>.420±.107</b>	<b>.641±.098</b>	<b>.848±.102</b>
Time Transfer	GAE	.384±.117	.619±.101	.828±.095
	GraphSAGE (unsp.)	.352±.121	.601±.105	.815±.093
	Graph Infomax	.369±.116	.606±.102	.821±.089
	GPT-GNN (Attr)	.374±.114	.614±.098	.826±.089
	GPT-GNN (Edge)	.397±.105	.629±.102	.836±.088
	GPT-GNN	<b>.405±.108</b>	<b>.635±.101</b>	<b>.840±.093</b>
Time + Field Transfer	GAE	.371±.124	.611±.108	.821±.102
	GraphSAGE (unsp.)	.349±.130	.602±.118	.812±.097
	Graph Infomax	.360±.121	.600±.102	.815±.093
	GPT-GNN (Attr)	.364±.115	.609±.103	.824±.094
	– (w/o node separation)	.347±.128	.601±.102	.813±.108
	GPT-GNN (Edge)	.390±.116	.622±.104	.830±.105
	– (w/o adaptive queue)	.376±.121	.617±.115	.828±.104
	GPT-GNN	<b>.397±.112</b>	<b>.628±.108</b>	<b>.833±.102</b>

- **All pre-training frameworks** help the performance of GNNs
  - GAE, GraphSage, Graph Infomax
  - GPT-GNN
- **GPT-GNN helps the most** by achieving a relative performance gain of 9.1% over the base model without pre-training
- **Both self-supervised tasks in GPT-GNN** help the pre-training framework
  - Attribute generation
  - Edge generation



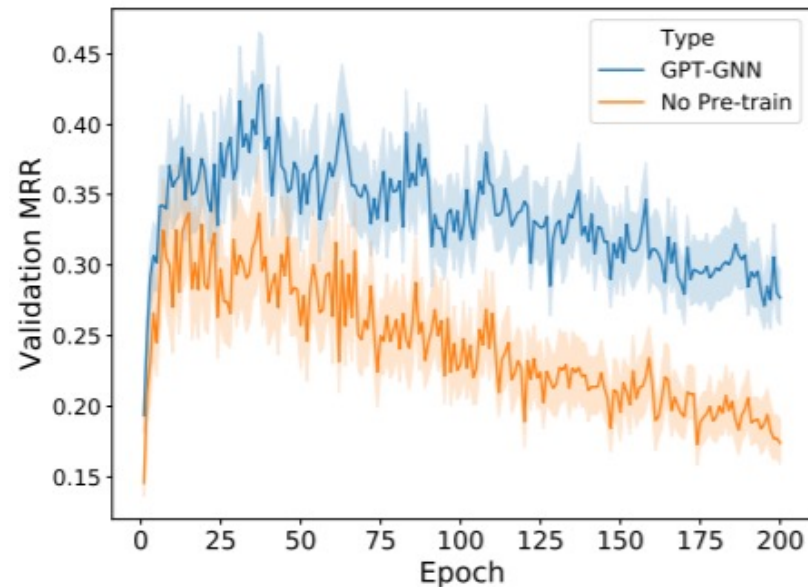
# GPT-GNN: Generative Pre-Training of GNNs

- Data: Microsoft Academic Graph

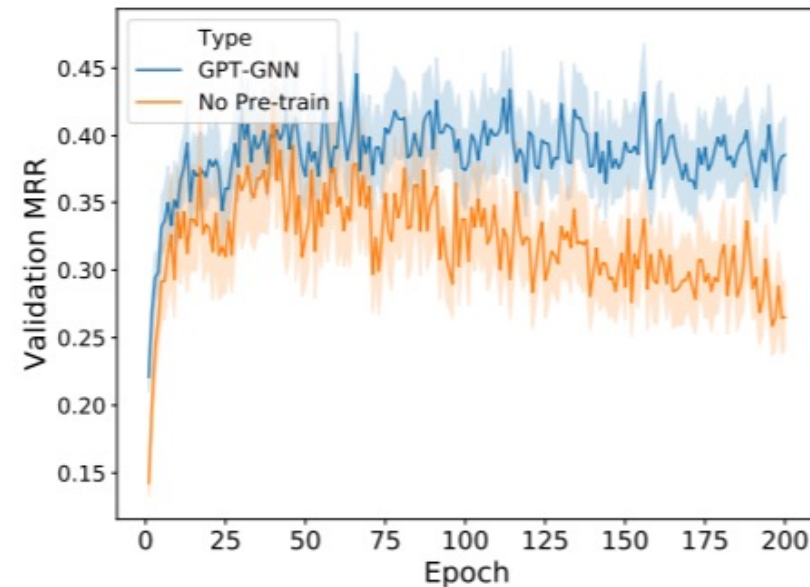


Model	HGT	GCN	GAT	RGCN	HAN
No Pre-train	<b>.346</b>	.327	.318	.296	.332
GPT-GNN	<b>.420</b>	.359	.382	.351	.406
Relative Gain	21.4%	9.8%	20.1%	18.9%	22.3%

# The Promise of Graph Pre-Training!



(a) Data Percentage: 10%



(b) Data Percentage: 20%

Predict Paper Title

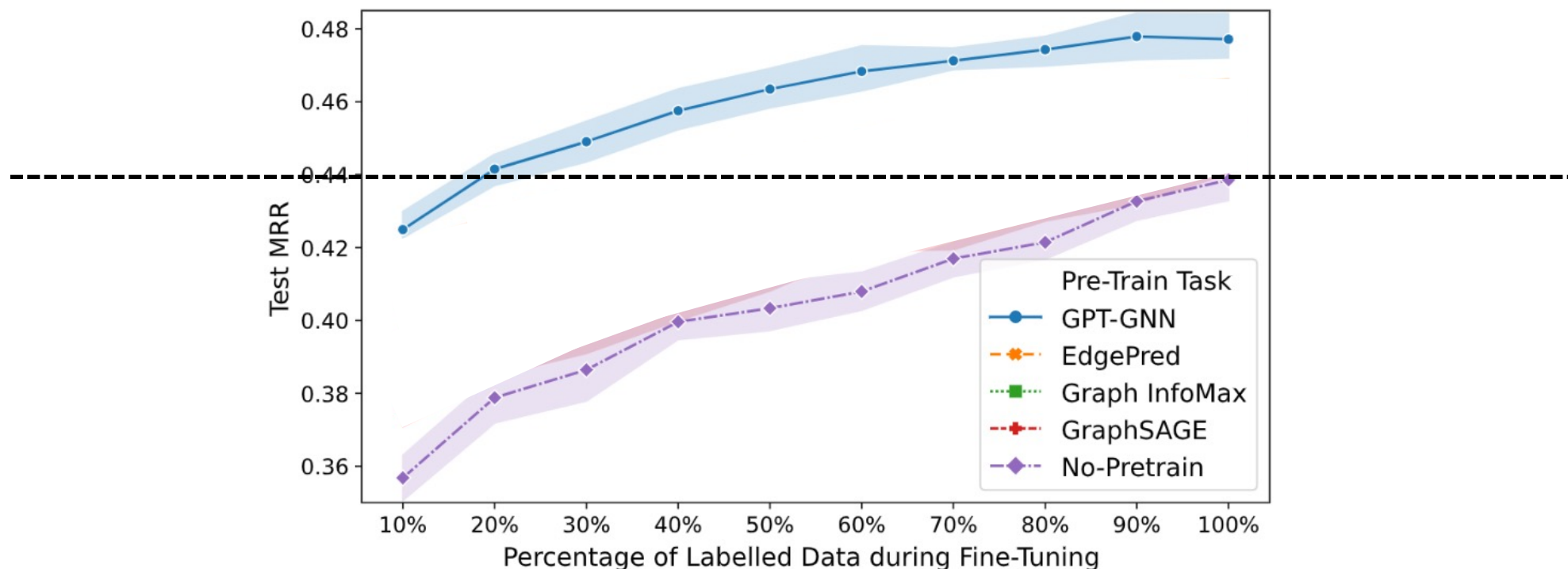
person recognition system using automatic probabilistic classification  
a novel framework using spectrum sensing in wireless systems  
a efficient evaluation of a distributed data storage service storage  
parameter control in wireless sensor networks networks networks  
a experimental system for for to the analysis of graphics

GroundTruth Paper Title

person re-identification by probabilistic relative distance comparison  
a secure collaborative spectrum sensing strategy in cyber physical systems  
an empirical analysis of a large scale mobile cloud storage service  
optimal parameter estimation under controlled communication over sensor networks  
an interactive computer graphics approach to surface representation

# The Promise of Graph Pre-Training!

During fine-tuning



The GNN model **w/o** pre-training with **100%** training data  
**VS**  
The pre-trained GNN model with **10-20%** training data

# Powering the Microsoft Office Graph



## One enterprise graph (monthly)

- 1.6 billion entities
  - 7 types of entities
- 7.8 trillion edges

# Anomaly detection on Microsoft Office Graph

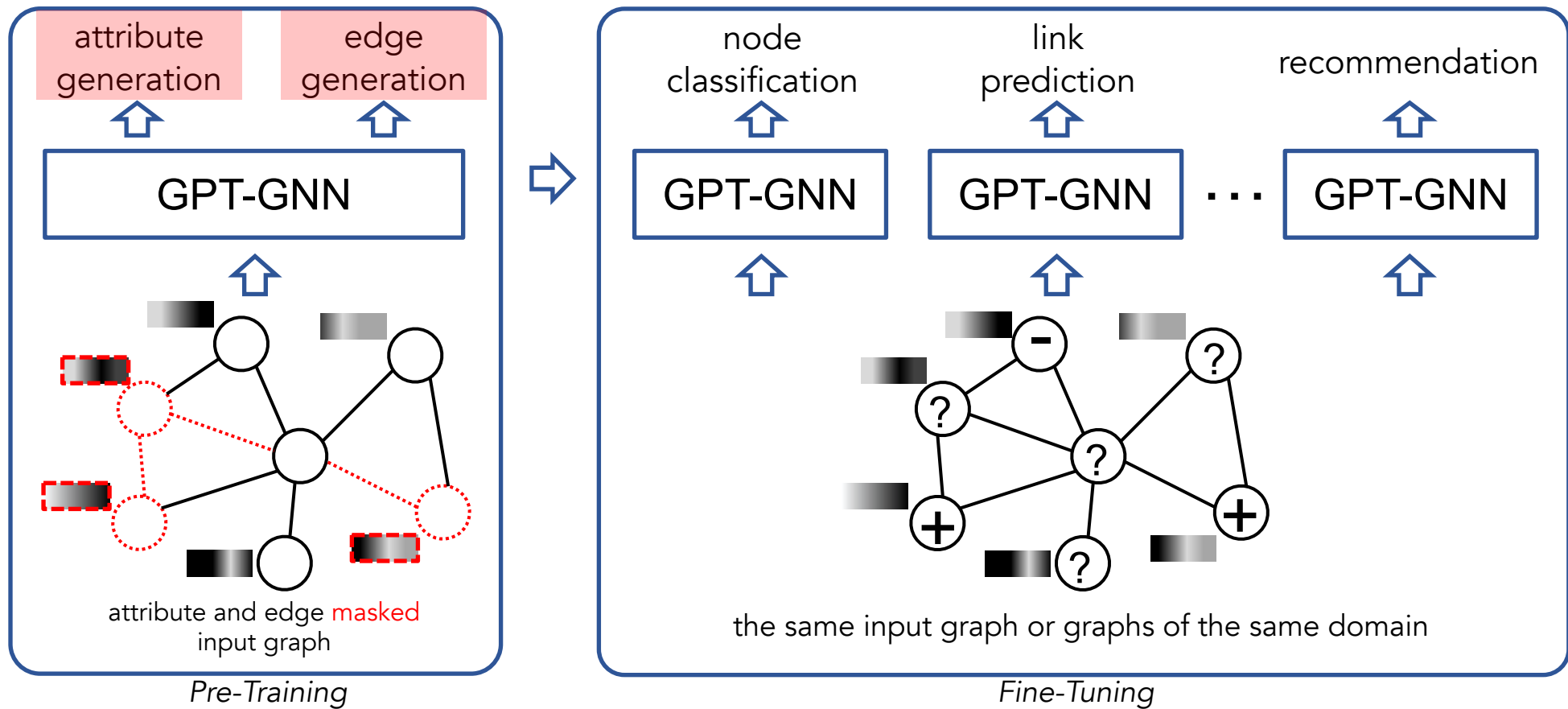
	Prec.	Recall	F1	Accu.
GraphSage	+0.00	+0.09	+0.06	+0.03
Graph Attention	+0.01	+0.11	+0.08	+0.03
HGT	+0.01	+0.30	+0.19	+0.07

Pre-trained  
HGT on  
one  
enterprise



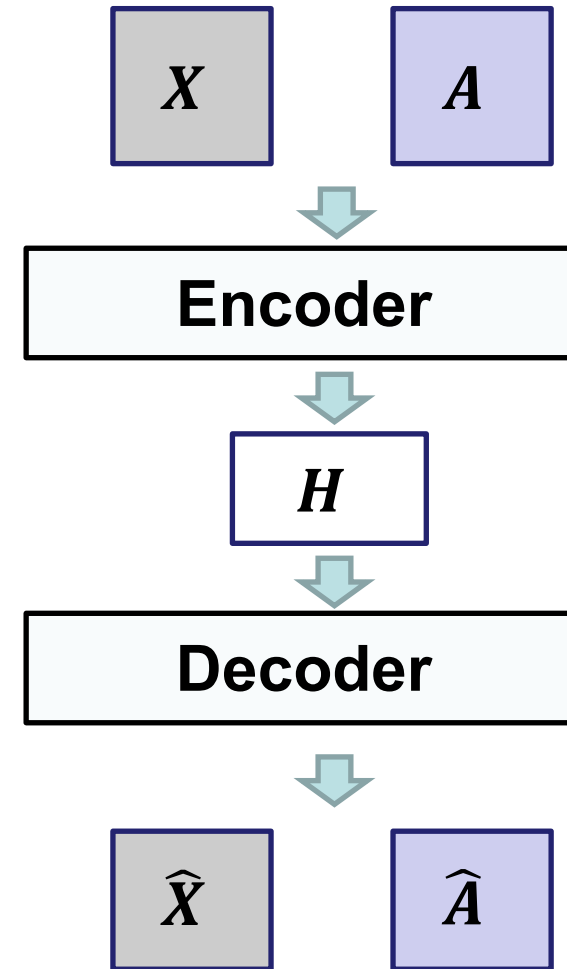
Other  
enterprise  
customers w/o  
data access

# GPT-GNN: Generative Pre-Training of GNNs



# Graph AutoEncoder

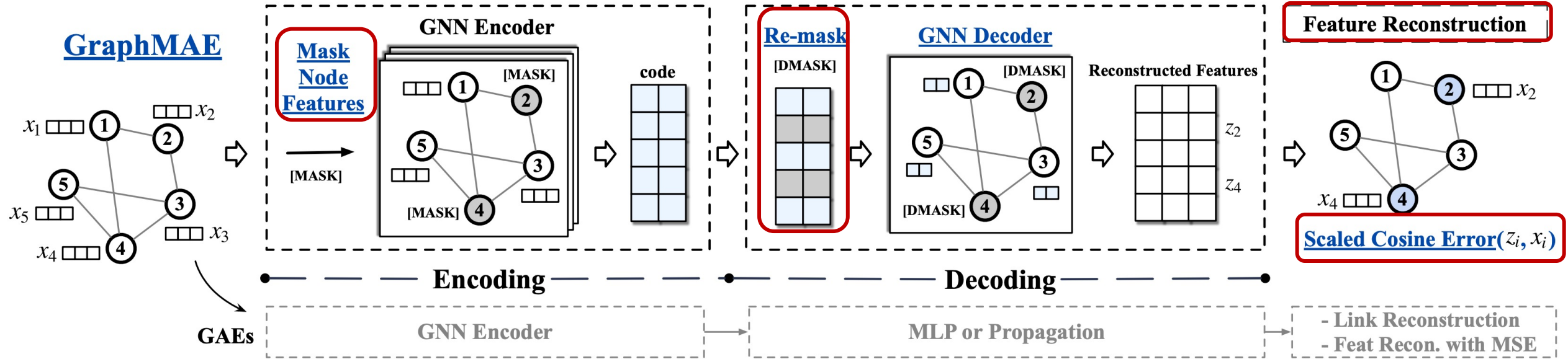
- $G = (V, A, X)$ 
  - $A \in \{0, 1\}^{N \times N}$ : adjacency matrix
  - $X \in \mathbb{R}^{N \times d}$ : node features
- Encoding
  - $H = f_E(A, X)$
- Decoding
  - $G' = f_D(A, H)$
- Reconstruction objectives
  - graph structure (link)
  - node features



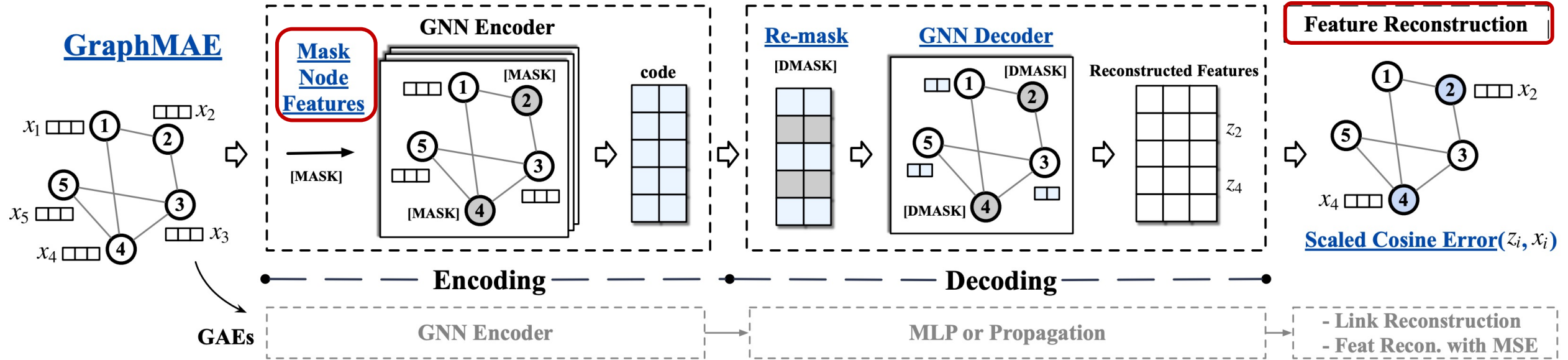
Methods	Reconstruction Target				Decoding Strategy		
	Feat. Loss	AE	No Struc.	Mask Feat.	GNN Decoder	Re-mask Dec.	Space
VGAE [20]	n/a	✓	-	-	-	-	$O(N^2)$
ARVGA [26]	n/a	✓	-	-	-	-	$O(N^2)$
MGAE [42]	MSE	✓	-	✓	-	-	$O(N)$
GALA [27]	MSE	✓	✓	-	✓	-	$O(N)$
GATE [31]	MSE	✓	-	-	✓	-	$O(N)$
AttrMask [16]	CE	✓	✓	✓	-	-	$O(N)$
GPT-GNN [17]	MSE	-	-	✓	-	-	$O(N)$
AGE [3]	n/a	✓	-	-	-	-	$O(N^2)$
NodeProp [18]	MSE	✓	✓	✓	-	-	$O(N)$
Error Function		Reconstruction Method					



# GraphMAE



# Masked Feature Reconstruction

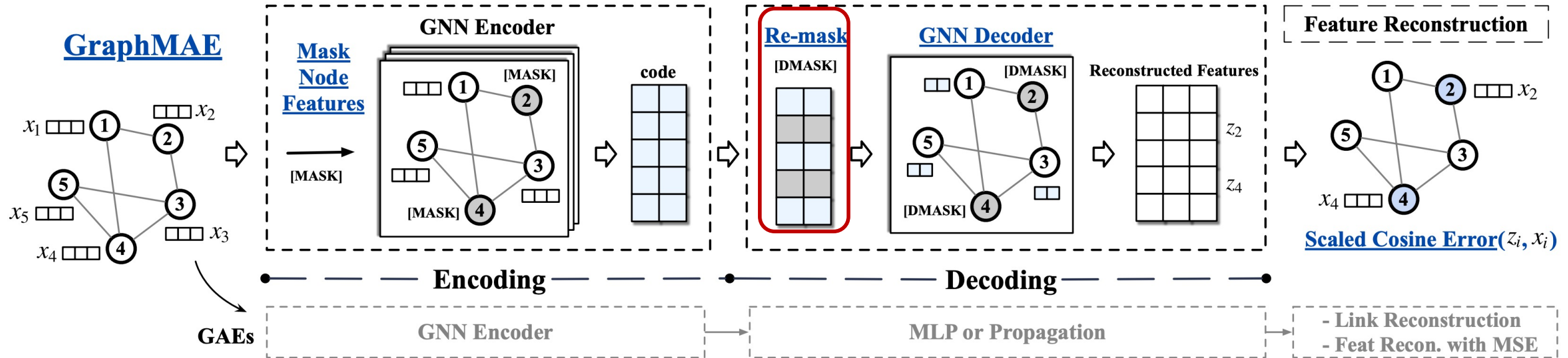


- Feature construction as the learning objective
- Masked feature reconstruction
  1. Sample a subset of nodes  $\tilde{V} \subset V$
  2. Replace node feature with [MASK]

$$\tilde{x}_i = \begin{cases} x_{[M]} & v_i \in \tilde{V} \\ x_i & v_i \notin \tilde{V} \end{cases}$$

- $H = J_E(A, \Lambda)$

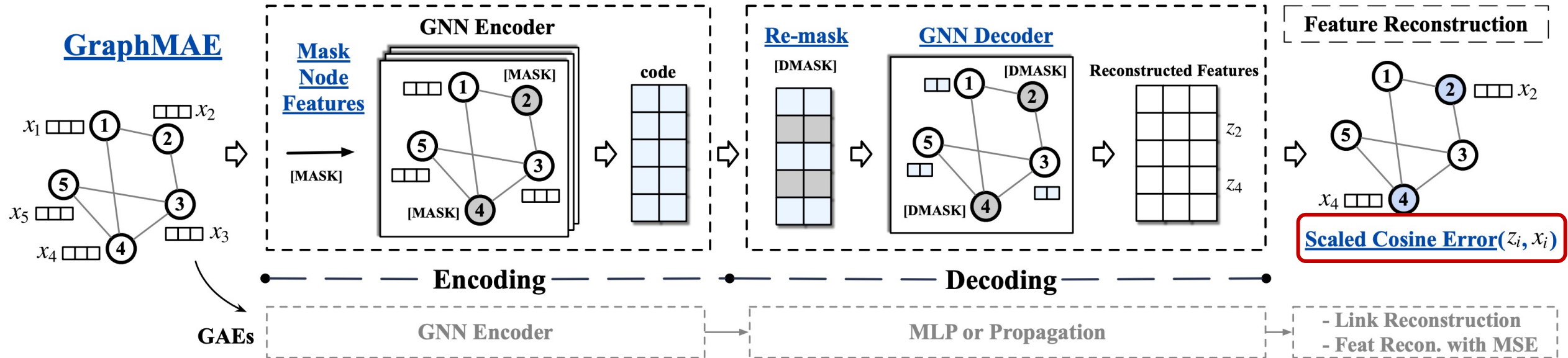
# GNNs as Decoder with Re-Mask Decoding



- Use a GNN as the decoder
  - A more expressive decoder helps reconstruct low informative features
- Re-mask node features before decoder
  - Re-mask the “masked” nodes

$$\tilde{H} = \text{Remask}(H), \quad Z = f_D(A, \tilde{H}) \quad \tilde{h}_i = \begin{cases} h_{[M]} & v_i \in \tilde{\mathcal{V}} \\ h_i & v_i \notin \tilde{\mathcal{V}} \end{cases}$$

# Scaled Cosine Error as the Criterion



- MSE fails, especially for continuous features
  - Sensitivity & low selectivity
- Scaled cosine error as the criterion
  - Cosine error & scaled coefficient

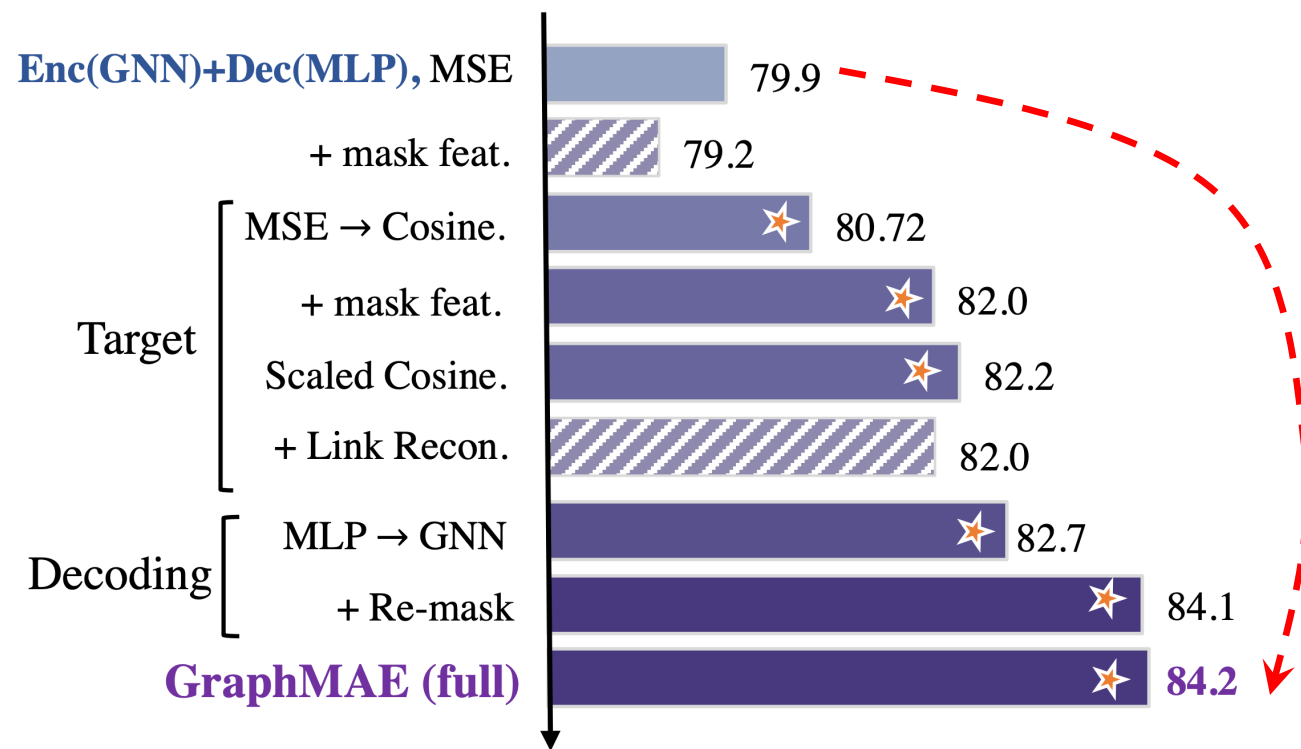
$$L_{MSE} = \frac{1}{|\tilde{V}|} \sum_{v_i \in \tilde{V}} (x_i - z_i)^2$$

$$\mathcal{L}_{SCE} = \frac{1}{|\tilde{\mathcal{V}}|} \sum_{v_i \in \tilde{\mathcal{V}}} \left(1 - \frac{x_i^T z_i}{\|x_i\| \cdot \|z_i\|}\right)^\gamma, \gamma \geq 1,$$

Methods	Reconstruction Target				Decoding Strategy		Space
	Feat. Loss	AE	No Struc.	Mask Feat.	GNN Decoder	Re-mask Dec.	
VGAE [20]	n/a	✓	-	-	-	-	$O(N^2)$
ARVGA [26]	n/a	✓	-	-	-	-	$O(N^2)$
MGAE [42]	MSE	✓	-	✓	-	-	$O(N)$
GALA [27]	MSE	✓	✓	-	✓	-	$O(N)$
GATE [31]	MSE	✓	-	-	✓	-	$O(N)$
AttrMask [16]	CE	✓	✓	✓	-	-	$O(N)$
GPT-GNN [17]	MSE	-	-	✓	-	-	$O(N)$
AGE [3]	n/a	✓	-	-	-	-	$O(N^2)$
NodeProp [18]	MSE	✓	✓	✓	-	-	$O(N)$
GraphMAE	SCE	✓	✓	✓	✓	✓	$O(N)$

Error Function      Reconstruction Method

# GraphMAE



(b) The effect of GraphMAE designs on the performance on Cora dataset.

**Table 4: Ablation studies of decoder type, re-mask and reconstruction criterion on node- and graph-level benchmarks.**

	Dataset	Node-Level			Graph-Level	
		Cora	PubMed	Arxiv	MUTAG	IMDB-B
COMP.	GraphMAE	84.2	81.1	71.75	88.19	75.52
	w/o mask	79.7	77.9	70.97	82.58	74.42
	w/o re-mask	82.7	80.0	71.61	86.29	74.42
	w/ MSE	79.1	73.1	67.44	86.30	74.04
Decoder	MLP	82.2	80.4	71.54	87.16	73.94
	GCN	81.3	79.1	71.59	87.78	74.54
	GIN	81.8	80.2	71.41	88.19	75.52
	GAT	84.2	81.1	71.75	86.27	74.04



# Downstream Tasks

## Node Classification

**Table 1: Experiment results in unsupervised representation learning for node classification.** We report Micro-F1(%) score for PPI and accuracy(%) for the other datasets.

	Dataset	Cora	CiteSeer	PubMed	Ogbn-arxiv	PPI	Reddit
Supervised	GCN	81.5	70.3	79.0	71.74±0.29	75.7±0.1	95.3±0.1
	GAT	83.0±0.7	72.5±0.7	79.0±0.3	72.10±0.13	97.30±0.20	96.0±0.1
Self-supervised	GAE	71.5±0.4	65.8±0.4	72.1±0.5	-	-	-
	GPT-GNN	80.1±1.0	68.4±1.6	76.3±0.8	-	-	-
	GATE	83.2±0.6	71.8±0.8	<u>80.9±0.3</u>	-	-	-
	DGI	82.3±0.6	71.8±0.7	76.8±0.6	70.34±0.16	63.80±0.20	94.0±0.10
	MVGRL	83.5±0.4	73.3±0.5	80.1±0.7	-	-	-
	GRACE <sup>1</sup>	81.9±0.4	71.2±0.5	80.6±0.4	71.51±0.11	69.71±0.17	94.72±0.04
	BGRL <sup>1</sup>	82.7±0.6	71.1±0.8	79.6±0.5	<u>71.64±0.12</u>	<u>73.63±0.16</u>	94.22±0.03
	InfoGCL	83.5±0.3	<b>73.5±0.4</b>	79.1±0.2	-	-	-
	CCA-SSG <sup>1</sup>	<u>84.0±0.4</u>	73.1±0.3	<u>81.0±0.4</u>	71.24±0.20	73.34±0.17	<u>95.07±0.02</u>
	GraphMAE	<b>84.2±0.4</b>	<u>73.4±0.4</u>	<b>81.1±0.4</b>	<b>71.75±0.17</b>	<b>74.50±0.29</b>	<b>96.01±0.08</b>

## Graph Classification

**Table 2: Experiment results in unsupervised representation learning for graph classification.** We report accuracy(%) for all datasets.

	Dataset	IMDB-B	IMDB-M	PROTEINS	COLLAB	MUTAG	REDDIT-B	NCI1
Supervised	GIN	75.1±5.1	52.3±2.8	76.2±2.8	80.2±1.9	89.4±5.6	92.4±2.5	82.7±1.7
	DiffPool	72.6±3.9	-	75.1±3.5	78.9±2.3	85.0±10.3	92.1±2.6	-
Graph Kernels	WL	72.30±3.44	46.95±0.46	72.92±0.56	-	80.72±3.00	68.82±0.41	80.31±0.46
	DGK	66.96±0.56	44.55±0.52	73.30±0.82	-	87.44±2.72	78.04±0.39	80.31±0.46
Self-supervised	graph2vec	71.10±0.54	50.44±0.87	73.30±2.05	-	83.15±9.25	75.78±1.03	73.22±1.81
	Infograph	73.03±0.87	49.69±0.53	74.44±0.31	70.65±1.13	89.01±1.13	82.50±1.42	76.20±1.06
	GraphCL	71.14±0.44	48.58±0.67	74.39±0.45	71.36±1.15	86.80±1.34	<u>89.53±0.84</u>	77.87±0.41
	JOAO	70.21±3.08	49.20±0.77	<u>74.55±0.41</u>	69.50±0.36	87.35±1.02	85.29±1.35	78.07±0.47
	GCC	72.0	49.4	-	78.9	-	<b>89.8</b>	-
	MVGRL	74.20±0.70	51.20±0.50	-	-	<u>89.70±1.10</u>	84.50±0.60	-
	InfoGCL	<u>75.10±0.90</u>	<u>51.40±0.80</u>	-	<u>80.00±1.30</u>	<b>91.20±1.30</b>	-	<u>80.20±0.60</u>
	GraphMAE	<b>75.52±0.66</b>	<b>51.63±0.52</b>	<b>75.30±0.39</b>	<b>80.32±0.46</b>	88.19±1.26	88.01±0.19	<b>80.40±0.30</b>

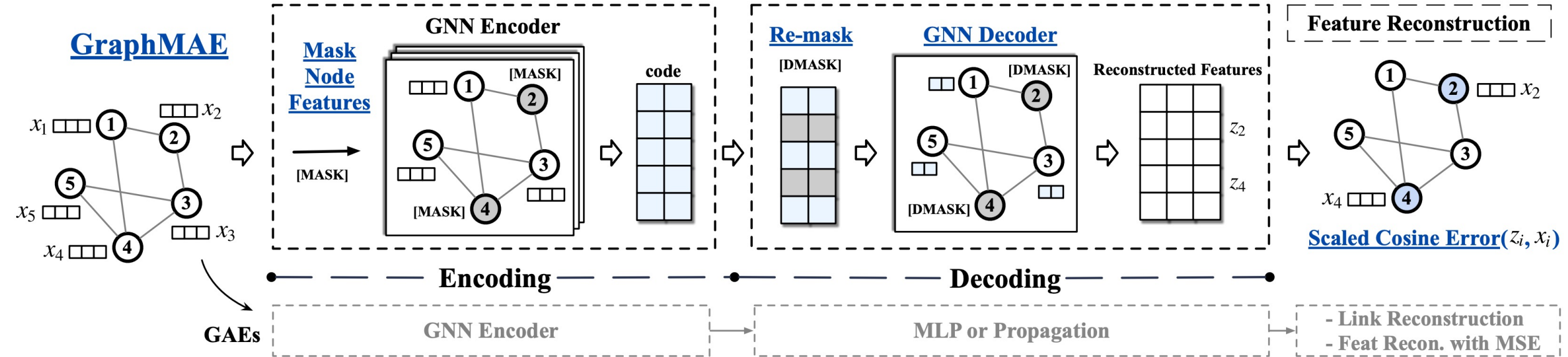
## Transfer Learning

**Table 3: Experiment results in transfer learning on molecular property prediction benchmarks.** The model is first pre-trained on ZINC15 and then finetuned on the following datasets. We report ROC-AUC(%) scores.

	BBBP	Tox21	ToxCast	SIDER	ClinTox	MUV	HIV	BACE	Avg.
No-pretrain	65.5±1.8	74.3±0.5	63.3±1.5	57.2±0.7	58.2±2.8	71.7±2.3	75.4±1.5	70.0±2.5	67.0
ContextPred	64.3±2.8	<u>75.7±0.7</u>	63.9±0.6	60.9±0.6	65.9±3.8	75.8±1.7	77.3±1.0	79.6±1.2	70.4
AttrMasking	64.3±2.8	<b>76.7±0.4</b>	<b>64.2±0.5</b>	<u>61.0±0.7</u>	71.8±4.1	74.7±1.4	77.2±1.1	79.3±1.6	71.1
Infomax	68.8 ±0.8	75.3 ±0.5	62.7 ±0.4	58.4 ±0.8	69.9±3.0	75.3 ±2.5	76.0 ±0.7	75.9 ±1.6	70.3
GraphCL	69.7±0.7	73.9±0.7	62.4±0.6	60.5±0.9	76.0±2.7	69.8±2.7	<b>78.5±1.2</b>	75.4±1.4	70.8
JOAO	70.2±1.0	75.0±0.3	62.9±0.5	60.0±0.8	<u>81.3±2.5</u>	71.7±1.4	76.7±1.2	77.3±0.5	71.9
GraphLoG	<b>72.5±0.8</b>	<u>75.7±0.5</u>	63.5±0.7	<b>61.2±1.1</b>	76.7±3.3	<u>76.0±1.1</u>	<u>77.8±0.8</u>	<b>83.5±1.2</b>	<u>73.4</u>
GraphMAE	<u>72.0±0.6</u>	75.5±0.6	<u>64.1±0.3</u>	60.3±1.1	<b>82.3±1.2</b>	<b>76.3±2.4</b>	77.2±1.0	<u>83.1±0.9</u>	<b>73.8</b>



# Summary of GraphMAE



1. Generative SSL on Graphs vs. Contrastive Learning on Graphs
2. Identify the common issues in current graph autoencoders
3. Present a simple masked graph autoencoder—**GraphMAE**

# Reflection & Motivation

- Problems in masked-feature-prediction
  - more *sensitive to the discriminability of input features*.

	Cora	PubMed
	raw $\rightarrow$ w/ PCA	raw $\rightarrow$ w/ PCA
Supervised	83.0 $\rightarrow$ 82.3 ( $\downarrow$ 0.7)	78.0 $\rightarrow$ 77.0 ( $\downarrow$ 1.0)
GraphMAE	84.2 $\rightarrow$ 82.6 ( $\downarrow$ 1.6)	81.1 $\rightarrow$ 78.9 ( $\downarrow$ 2.2)
GraphMAE2	84.5 $\rightarrow$ 83.5 ( $\downarrow$ 1.0)	81.4 $\rightarrow$ 80.1 ( $\downarrow$ 1.3)

- *raw* : the original node features
- *w/ PCA* : the input features are reduced to 50-dimensional vectors using PCA

**Resolution:** imposing regularization on target reconstruction

# GraphMAE2

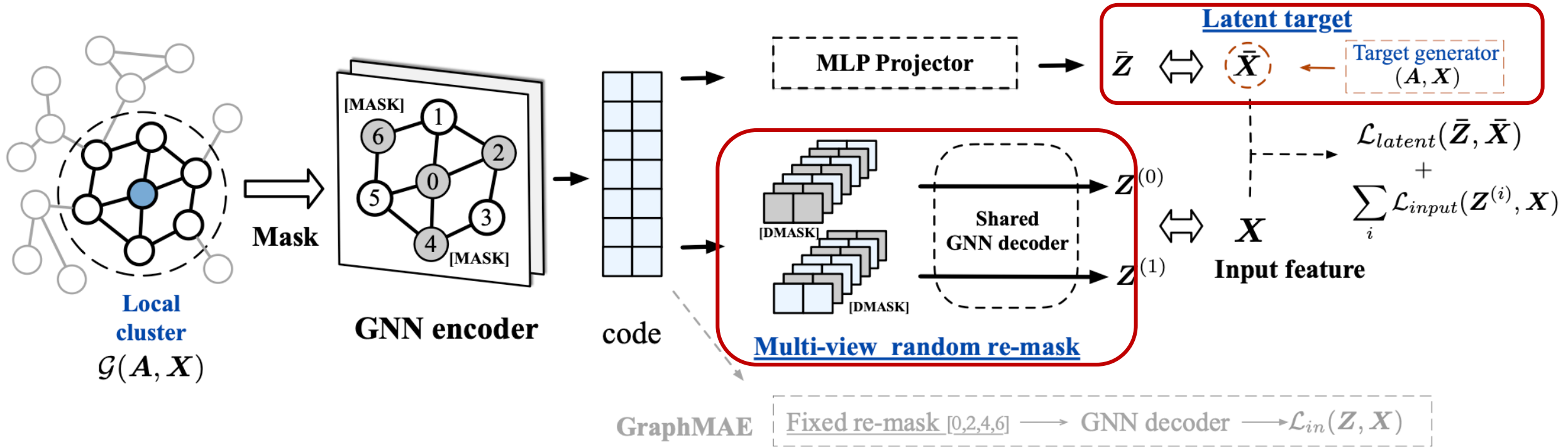
A Decoding-Enhanced Masked Self-Supervised Graph Learner.

10:00-10:10 AM

Thursday, May 4, 2023

@Classroom 107

# GraphMAE2 Framework



- Multi-view random re-mask decoding
- Latent representation prediction
- Scaling to large-scale graphs with local clustering

# Multi-View Random Re-Mask Decoding

– Avoid representation overfitting to input features

– **Randomly** re-mask representations/code

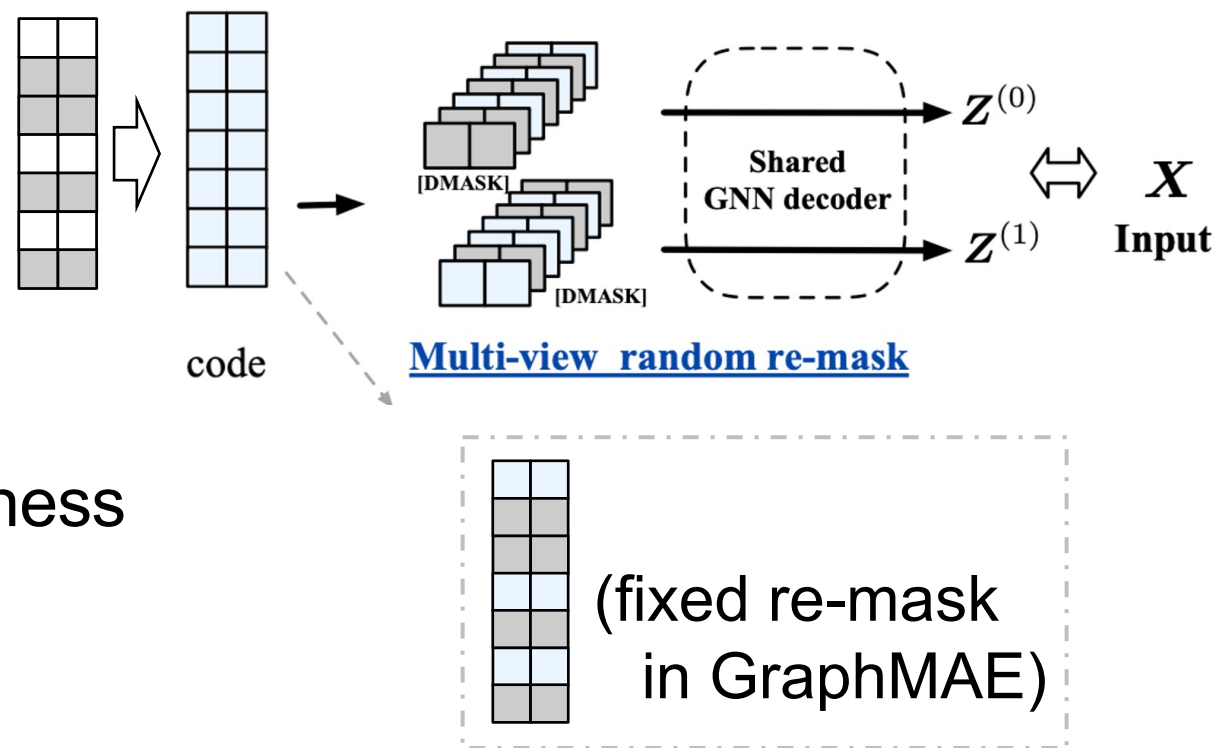
- $\tilde{H} = \text{Remask}(H), Z = f_D(A, \tilde{H})$

$$\tilde{h}_i = \begin{cases} \mathbf{h}_{[M]} & v_i \in \overline{\mathcal{V}} \\ \mathbf{h}_i & v_i \notin \overline{\mathcal{V}} \end{cases}$$

– **Multiple** re-masking

- $K$ -different randomly re-masking
- better generalization and effectiveness

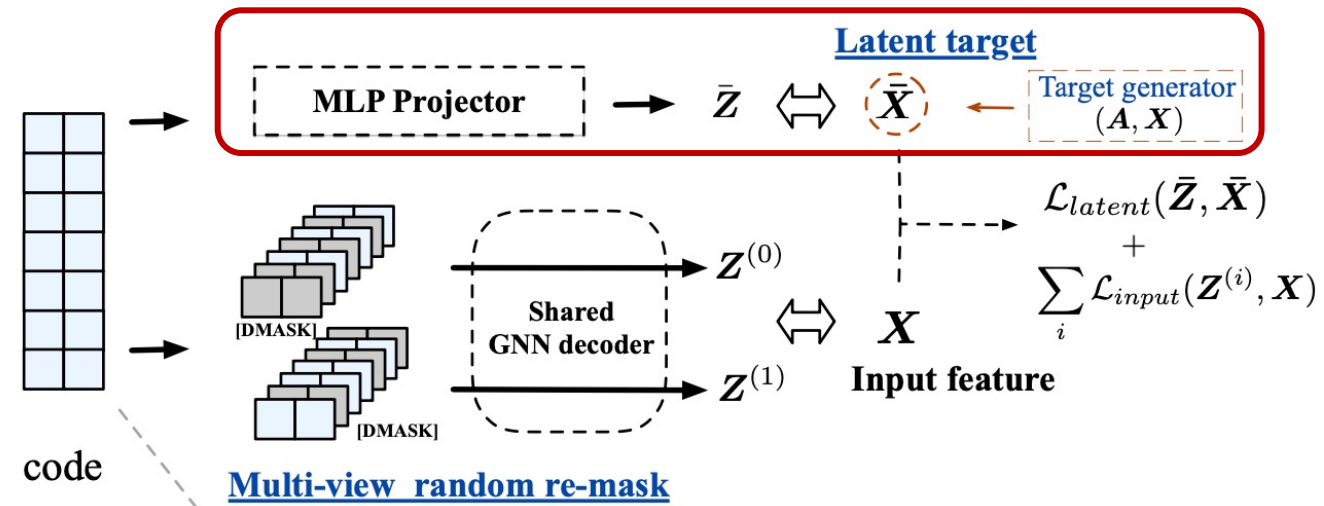
$$\mathcal{L}_{input} = \frac{1}{|\tilde{\mathcal{V}}|} \sum_{j=1}^K \sum_{v_i \in \tilde{\mathcal{V}}} \left(1 - \frac{\mathbf{x}_i^\top \mathbf{z}_i^{(j)}}{\|\mathbf{x}_i\| \cdot \|\mathbf{z}_i^{(j)}\|}\right)^Y$$



# Latent Representation Prediction

- Additional informative prediction target
  - Minimally affected by input features & and GNN as a denoiser
- Predicting **masked latent representations**
  - A (*momentum*) target generator  $f_{target}(\cdot|\xi)$
  - Prediction:  $\bar{Z} = f_E(\text{mask}(G); \theta)$
  - Latent target:  $\bar{X} = f_{target}(G; \xi)$ 
    - $\xi \leftarrow \tau \cdot \xi + (1 - \tau) \cdot \theta$

$$\mathcal{L}_{latent} = \frac{1}{N} \sum_i \left(1 - \frac{\bar{z}_i^\top \bar{x}_i}{\|\bar{z}\| \cdot \|\bar{x}\|}\right)^\gamma$$



# Experiments

- Embedding evaluation with linear probing
- On large-scale graphs

Datasets	#Nodes	#Edges
Cora	2,485	5,069
Citeseer	2,110	3,668
Pubmed	19,717	44,324
ogbn-Arxiv	169,343	1,166,243
ogbn-Products	2,449,029	61,859,140
MAG-Scholar-F	12,403,930	358,010,024
ogbn-Papers100M	111,059,956	1,615,685,872



# Linear Probing

- **Setting: training a linear classifier**
- Results:
  1. GraphMAE2 consistently outperforms all baselines
  2. *Random-Init* models can achieve decent results

	Arxiv	Products	MAG	Papers100N
MLP	55.50±0.23	61.06±0.08	39.11±0.21	47.24±0.31
SGC	66.92±0.08	74.87±0.25	54.68±0.23	63.29±0.19
Random-Init	68.14±0.02	74.04±0.06	56.57±0.03	61.55±0.12
CCA-SSG	68.57±0.02	75.27±0.05	51.55±0.03	55.67±0.15
GRACE	69.34±0.01	<u>79.47±0.59</u>	57.39±0.02	61.21±0.12
BGRL	70.51±0.03	78.59±0.02	57.57±0.01	62.18±0.15
GGD <sup>1</sup>	-	75.70±0.40	-	<u>63.50±0.50</u>
GraphMAE	<u>71.03±0.02</u>	78.89±0.01	<u>58.75±0.03</u>	62.54±0.09
GraphMAE2	<b>71.89±0.03</b>	<b>81.59±0.02</b>	<b>59.24±0.01</b>	<b>64.89±0.04</b>

	Cora	CiteSeer	PubMed
GCN	81.5	70.3	79.0
GAT	83.0±0.7	72.5±0.7	79.0±0.3
GAE	71.5±0.4	65.8±0.4	72.1±0.5
DGI	82.3±0.6	71.8±0.7	76.8±0.6
MVGRL	83.5±0.4	73.3±0.5	80.1±0.7
GRACE	81.9±0.4	71.2±0.5	80.6±0.4
BGRL	82.7±0.6	71.1±0.8	79.6±0.5
InfoGCL	83.5±0.3	<b>73.5±0.4</b>	79.1±0.2
CCA-SSG	84.0±0.4	73.1±0.3	81.0±0.4
GGD	83.9±0.4	73.0±0.6	<u>81.3±0.8</u>
GraphMAE	<u>84.2±0.4</u>	<u>73.4±0.4</u>	81.1±0.4
GraphMAE2	<b>84.5±0.6</b>	<u>73.4±0.3</u>	<b>81.4±0.5</b>

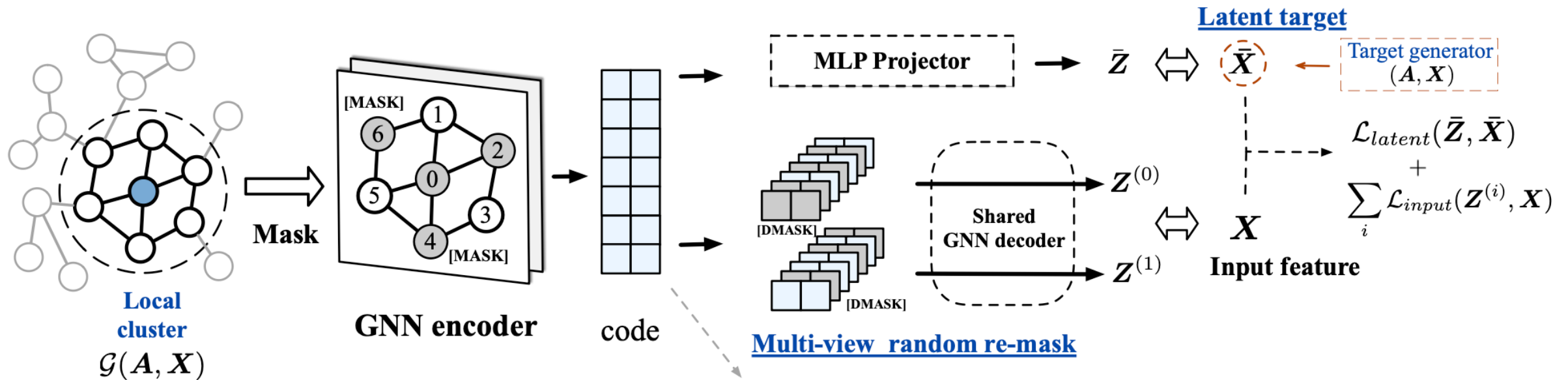
# Ablation Studies

## Component Ablation of GraphMAE2

- GraphMAE2 surpasses all baselines with the same sampling strategy
- Using local clusters brings further improvement

	Strategy	Products	MAG	Papers100M
GRACE	<i>SAINT</i>	79.47 $\pm$ 0.59	57.39 $\pm$ 0.02	61.21 $\pm$ 0.12
BGRL	<i>SAINT</i>	78.59 $\pm$ 0.02	57.57 $\pm$ 0.01	62.18 $\pm$ 0.15
GraphMAE2	<i>SAINT</i>	80.96 $\pm$ 0.03	58.75 $\pm$ 0.03	64.21 $\pm$ 0.11
GraphMAE2	<i>Cluster</i>	79.35 $\pm$ 0.05	58.05 $\pm$ 0.02	63.77 $\pm$ 0.11
GraphMAE2	<i>LC</i>	81.59 $\pm$ 0.02	59.24 $\pm$ 0.01	64.89 $\pm$ 0.12

# GraphMAE2 Summary

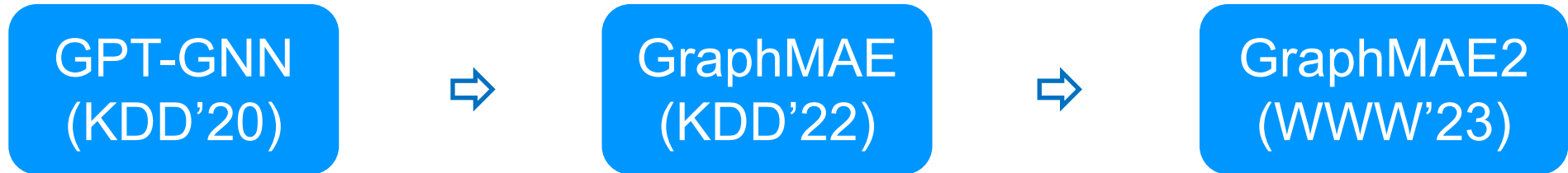


- Analyze the problem in masked feature prediction
- Present GraphMAE2 with improved decoding strategies
- GraphMAE2 achieves promising performance in large-scale graphs



GraphMAE2: <https://github.com/THUDM/GraphMAE2>  
GraphMAE: <https://github.com/THUDM/GraphMAE>

# Generative Learning on Graphs



# Pre-Train Graphs with *Language/Image/Knowledge*



Academic Graph



Microsoft Office Graph

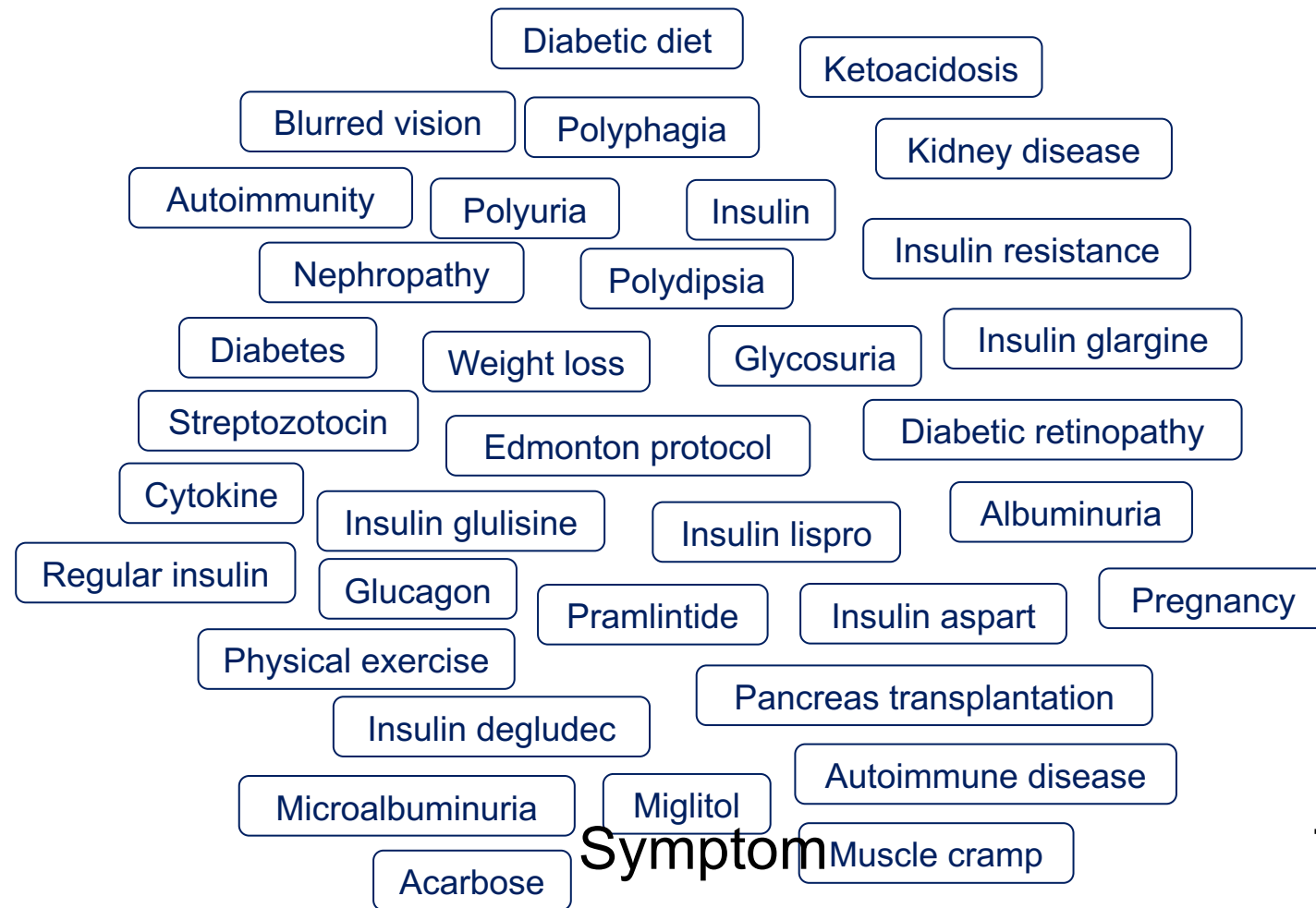


LinkedIn Economic Graph



Facebook Entity Graph

# Neural Symbolic Reasoning



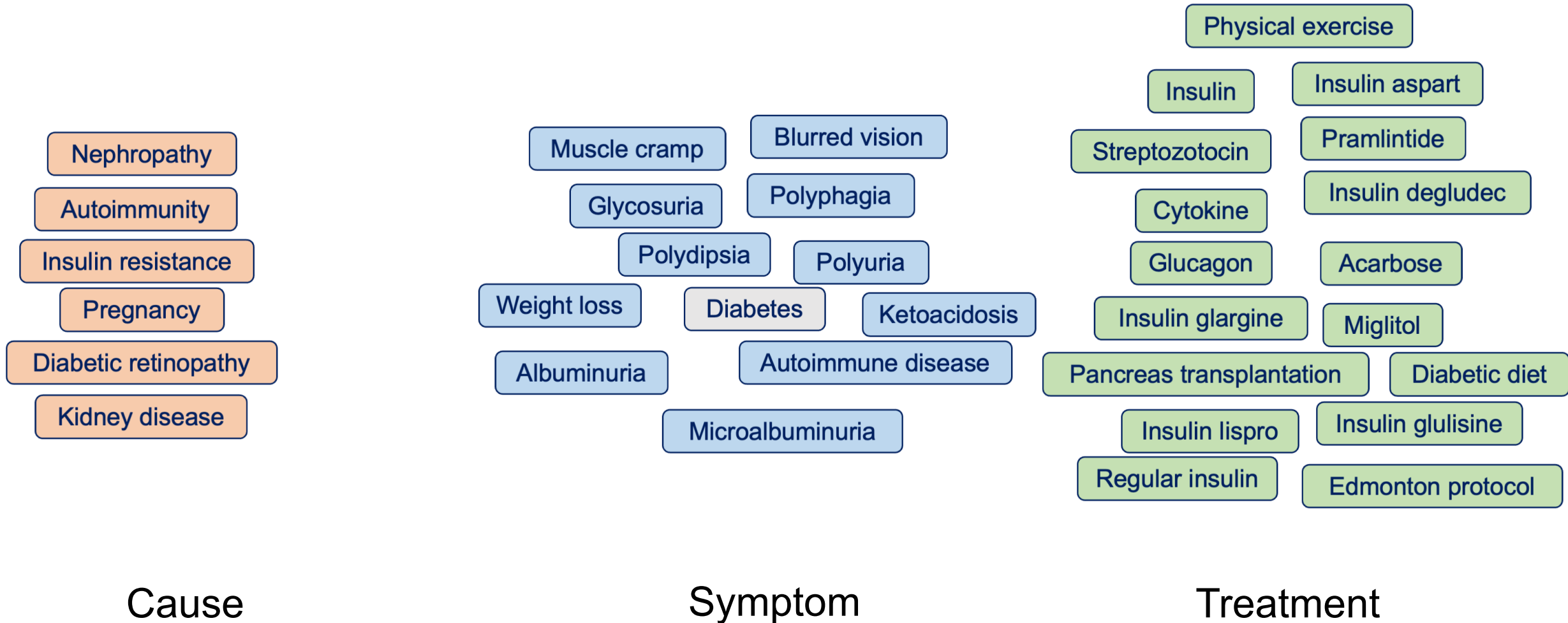
Cause

Symptom

Treatment

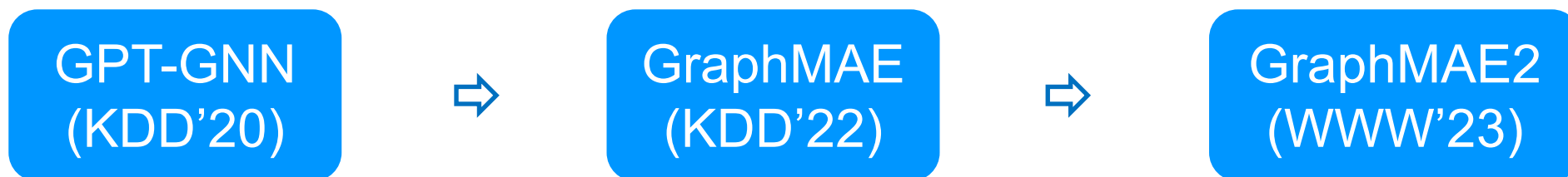


# Neural Symbolic Reasoning



# GNNs vs. LLMs

GNN:



LLM:

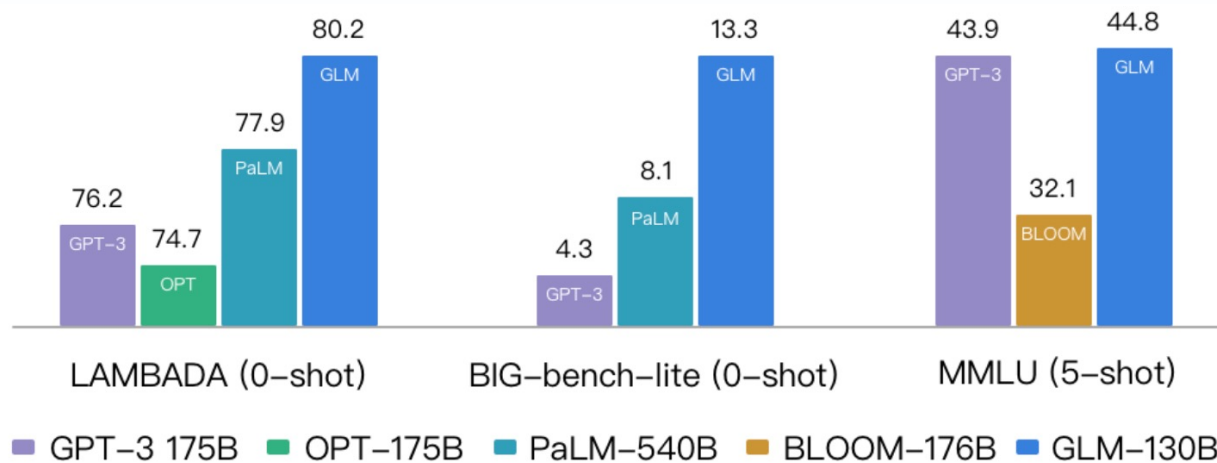


<https://github.com/THUDM>

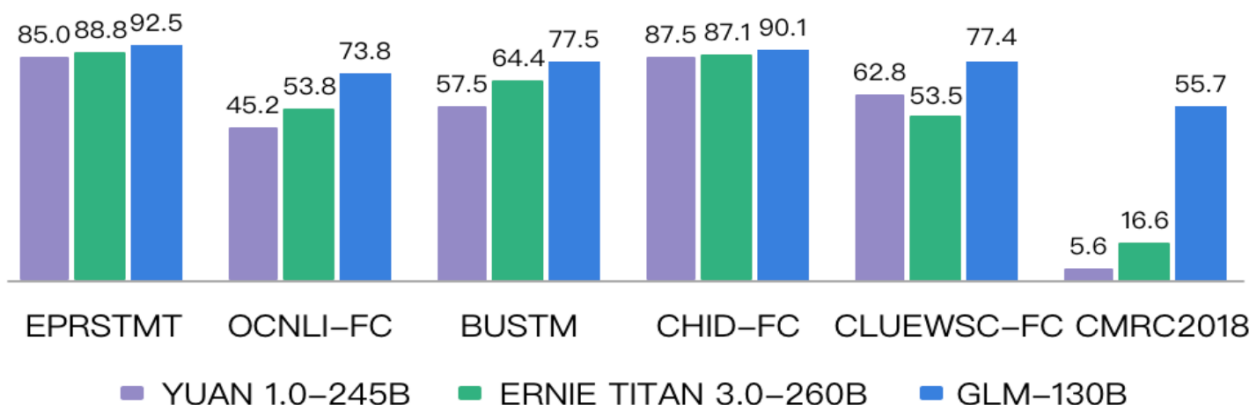


# GLM-130B

## English: better than GPT-3/OPT/PaLM on MMLU, LAMBADA, BIG-bench-lite



## Chinese: better than ERNIE 260B & YUAN 245B



Since Aug., 2022, requests from ~1000 orgs in 69 countries

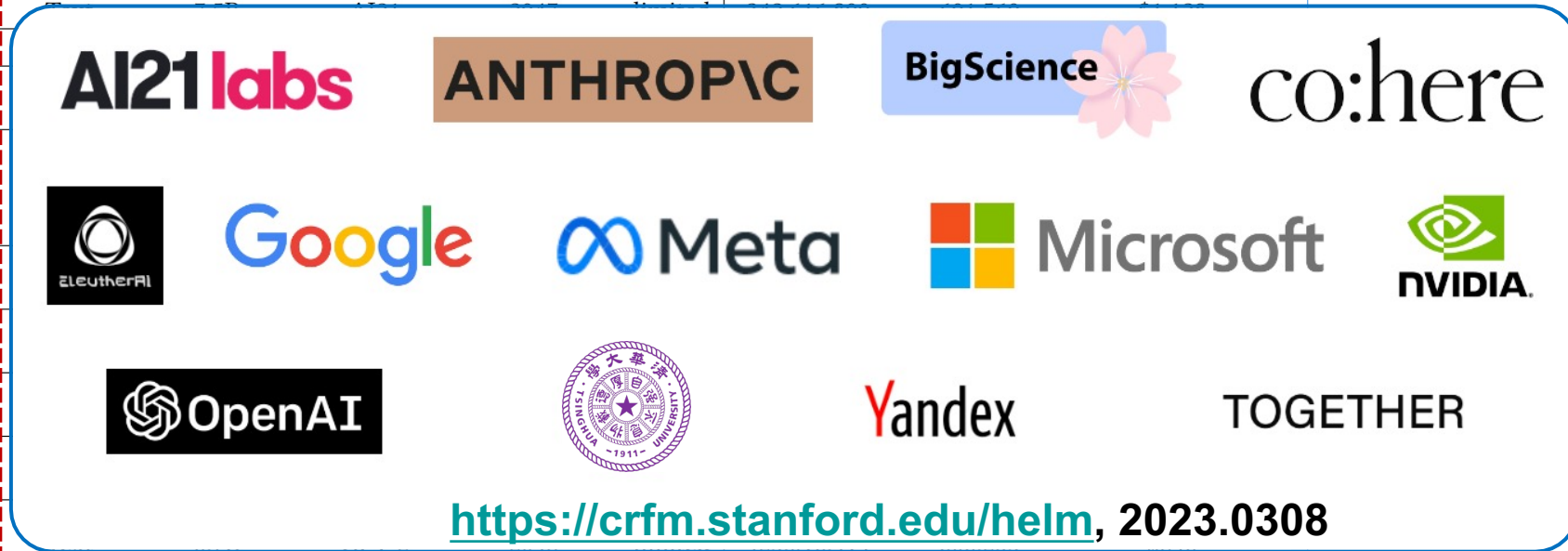
- Google
- Microsoft
- Facebook
- Stanford
- MIT
- UC Berkely
- CMU
- Harvard
- Princeton
- Yale
- Cornell
- UIUC
- Cambridge
- Oxford
- Huawei
- Alibaba
- Tencent
- Baidu
- Meituan
- Bytedance
- Didi
- Xiaoice
- Xiaodu
- Xiaomi
- Xiaopeng
- Youdao
- Face++
- Ping An Cap
- Peking U.
- Zhejiang U.
- Shanghai JT U.
- Fudan U.
- USTC
- U of CAS
- Wuhan U.
- Naikai U.
- Hongkong U.
- CUHK
- HKUST
- BAAI
- Zhejiang Lab
- Shanghai AI Lab

# GLM-130B

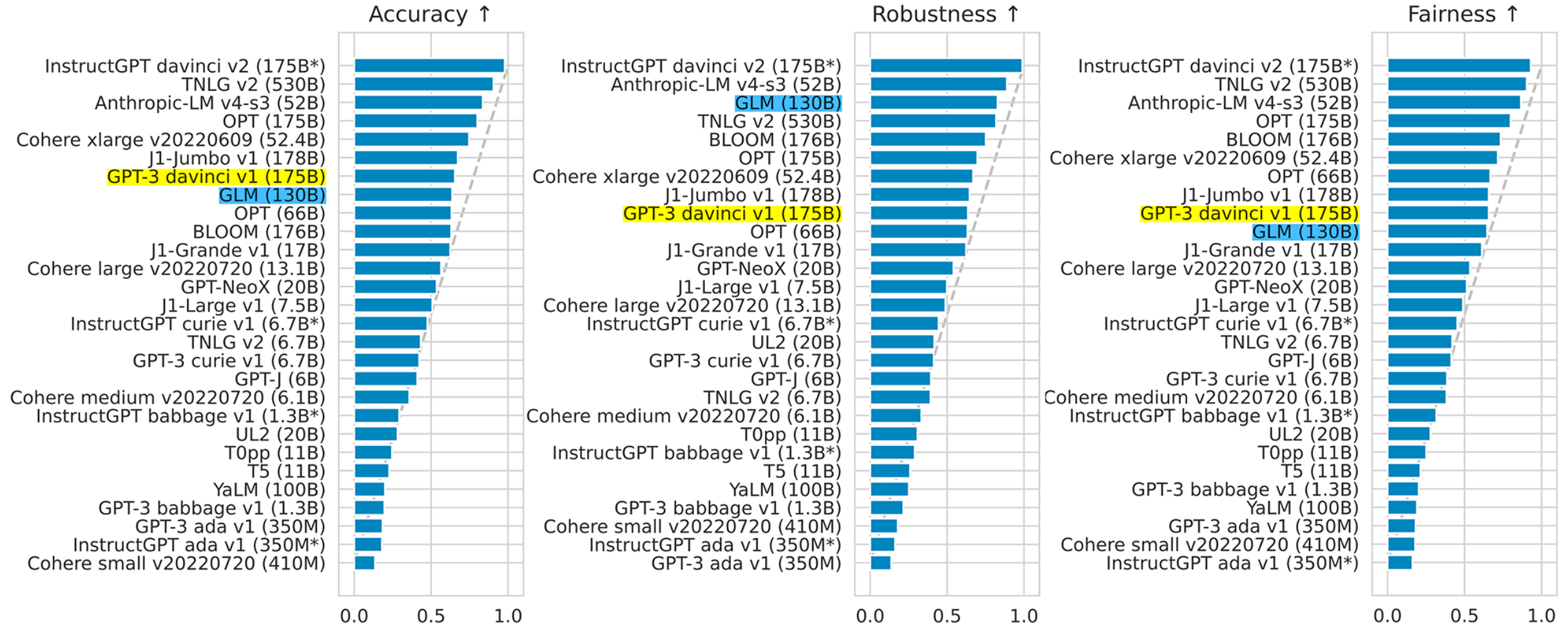


The only model from academia was covered by Stanford's HELM

Model	Model Creator	Modality	# Parameters	Tokenizer	Window Size	Access	Total Tokens	Total Queries	Total Cost
J1-Jumbo v1 (178B)	AI21 Labs	Text	178B	AI21	2047	limited	327,443,515	591,384	\$10,926
J1-Grande v1 (17B)	AI21 Labs	Text	17B	AI21	2047	limited	326,815,150	591,384	\$2,973
J1-Large v1 (7.5B)	AI21 Labs								
Anthropic-LM v4-s3 (52B)	Anthropic								
BLOOM (176B)	BigScience								
T0++ (11B)	BigScience								
Cohere xlarge v20220609 (52.4B)	Cohere								
Cohere large v20220720 (13.1B) <sup>58</sup>	Cohere								
Cohere medium v20220720 (6.1B)	Cohere								
Cohere small v20220720 (410M) <sup>59</sup>	Cohere								
GPT-J (6B)	EleutherAI								
GPT-NeoX (20B)	EleutherAI								
T5 (11B)	Google								
UL2 (20B)	Google								
OPT (66B)	Meta								
OPT (175B)	Meta								
TNLG v2 (6.7B)	Microsoft/NVIDIA								
TNLG v2 (530B)	Microsoft/NVIDIA								
GPT-3 davinci v1 (175B)	OpenAI								
GPT-3 curie v1 (6.7B)	OpenAI								
GPT-3 babbage v1 (1.3B)	OpenAI	Text	1.3B	GPT-2	2048	limited	422,123,900	606,253	\$211
GPT-3 ada v1 (350M)	OpenAI	Text	350M	GPT-2	2048	limited	422,635,705	604,253	\$169
InstructGPT davinci v2 (175B*)	OpenAI	Text	175B*	GPT-2	4000	limited	466,872,228	599,815	\$9,337
InstructGPT curie v1 (6.7B*)	OpenAI	Text	6.7B*	GPT-2	2048	limited	420,004,477	606,253	\$840
InstructGPT babbage v1 (1.3B*)	OpenAI	Text	1.3B*	GPT-2	2048	limited	419,036,038	604,253	\$210
InstructGPT ada v1 (350M*)	OpenAI	Text	350M*	GPT-2	2048	limited	418,915,281	604,253	\$168
Codex davinci v2	OpenAI	Code	Unknown	GPT-2	4000	limited	46,272,590	57,051	\$925
Codex cushman v1	OpenAI	Code	Unknown	GPT-2	2048	limited	42,659,399	59,751	\$85
GLM (130B)	Tsinghua University	Text	130B	ICE	2048	open	375,474,243	406,072	2,100 GPU hours
YaLM (100B)	Yandex	Text	100B	Yandex	2048	open	378,607,292	405,093	2,200 GPU hours



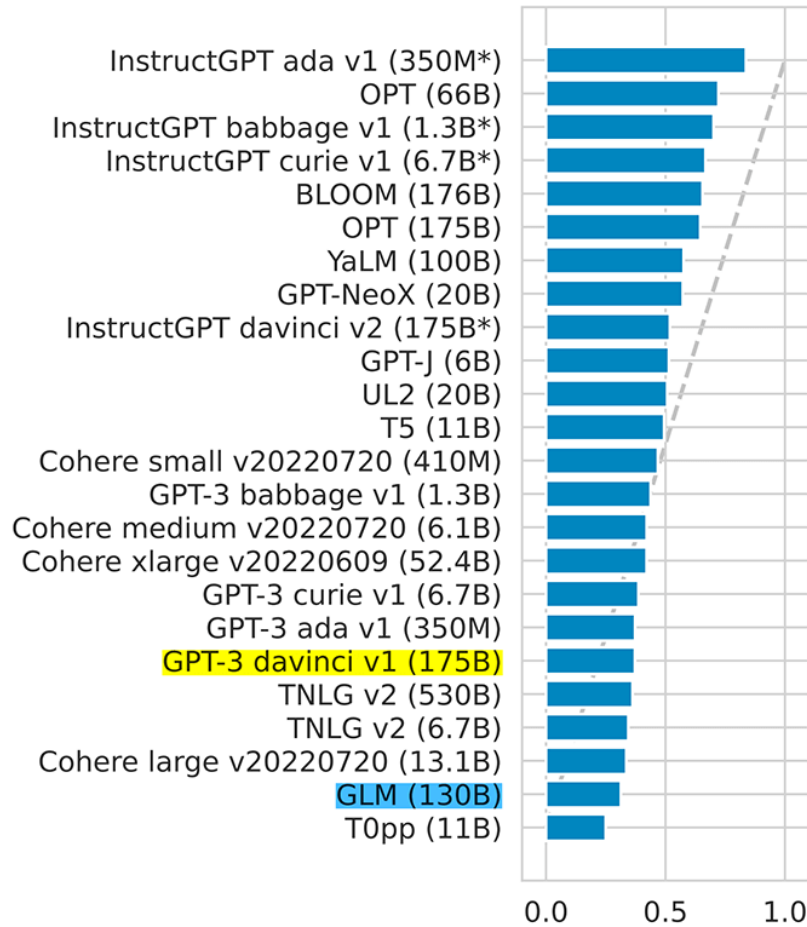
# GLM-130B



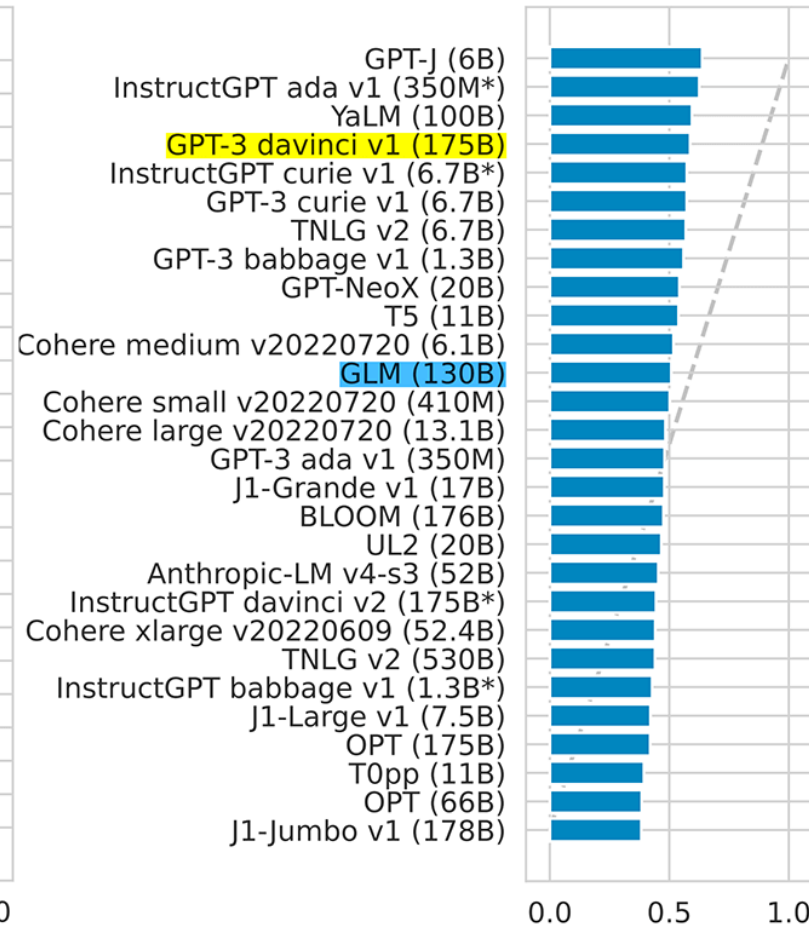
# GLM-130B



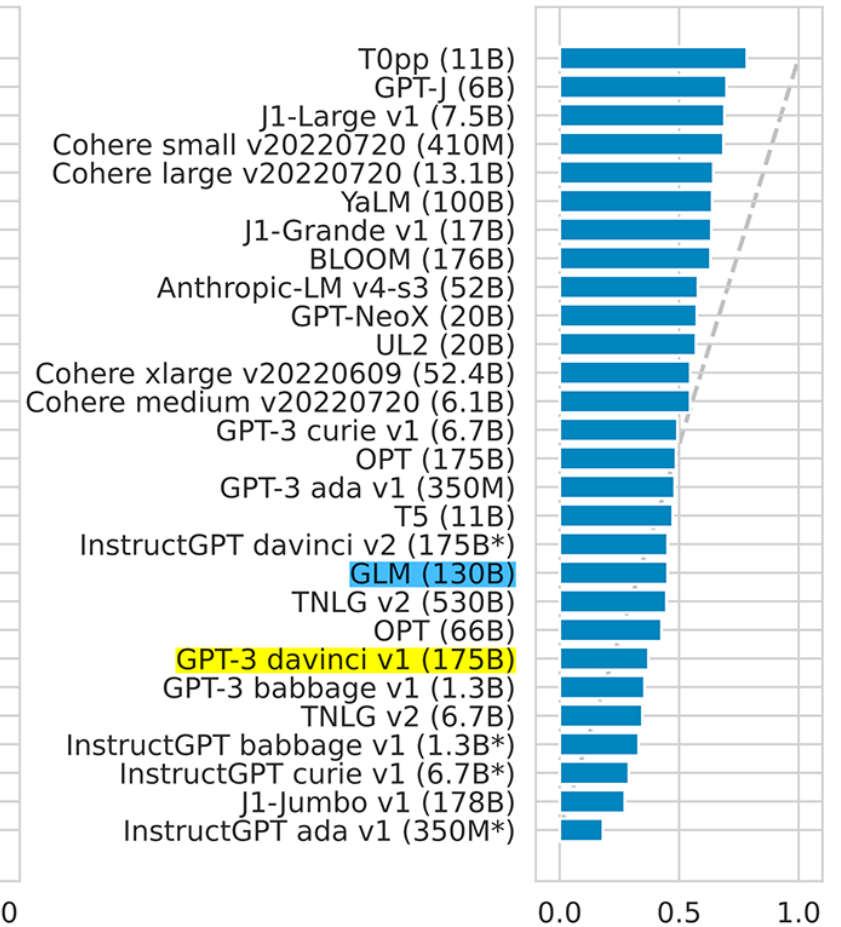
Calibration error ↓



Bias ↓



Toxicity ↓





# ChatGLM

Alpha 内测

提出你的想法



🕒 收养一只流浪猫好还是买一只...

🕒 2022年世界杯冠军是谁

...

## 💡 试试这些例子：

列出一些年夜饭好意头的菜肴以及其寓意。

帮我写一篇人工智能课程的教案，1000字。

怎么修改huggingface transformers的model cache位置？

## ⚠️ 当前模型限制：

可能会生成不正确的信息

可能会产生有害说明或有偏见的内容

暂时不擅长逻辑类回答，如数学和编程类问题

Alpha/beta test from Mar. 14, 2023


# An Open ChatGLM-6B




- Mar. 14, 2023, open-sourced model
- Mar. 16, 2023, **#1** on GitHub Trending
- Mar. 18-30, **#1** on Hugging Face Trending
- Apr. 30, 2023, 21.8k stars in GitHub

**1M** downloads in HF

 THUDM/chatglm-6b  
Updated 3 days ago • ↓ 910k • ❤️ 1.76k

 THUDM/chatglm-6b-int4  
Updated 3 days ago • ↓ 91.9k • ❤️ 199

 THUDM/chatglm-6b-int8  
Updated 3 days ago • ↓ 14.1k • ❤️ 12



<https://huggingface.co/THUDM>

# Thank you !



**ChatGLM-6B** Public



ChatGLM-6B: An Open Bilingual Dialogue Language Model | 开源双语对话语言模型

Python 21.9k 2.6k



**GLM-130B** Public



GLM-130B: An Open Bilingual Pre-Trained Model (ICLR 2023)

Python 5.1k 365



**CodeGeeX** Public



CodeGeeX: An Open Multilingual Code Generation Model

Python 4.8k 316



**CogView** Public



Text-to-Image generation. The repo for NeurIPS 2021 paper "CogView: Mastering Text-to-Image Generation via Transformers".

Python 1.4k 163



**CogVideo** Public



Text-to-video generation. The repo for ICLR2023 paper "CogVideo: Large-scale Pretraining for Text-to-Video Generation via Transformers"

Python 2.8k 286



**cogdl** Public



CogDL: A Comprehensive Library for Graph Deep Learning (WWW 2023)

Python 1.4k 300



slides



<https://github.com/THUDM>



@thukeyg