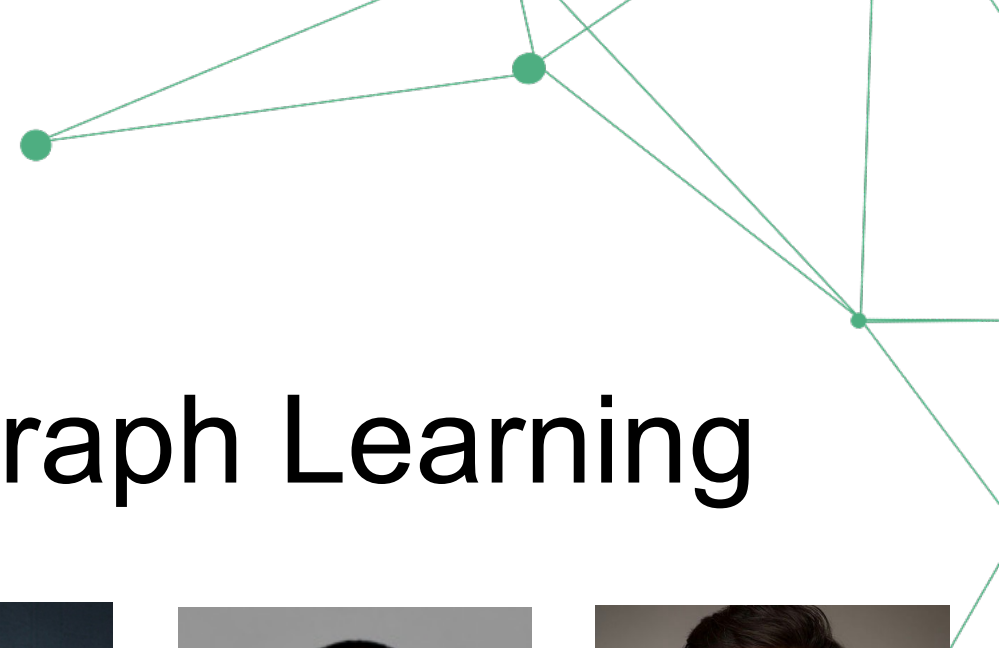




AAAI-24 Tutorial  
Feb 2024, Vancouver



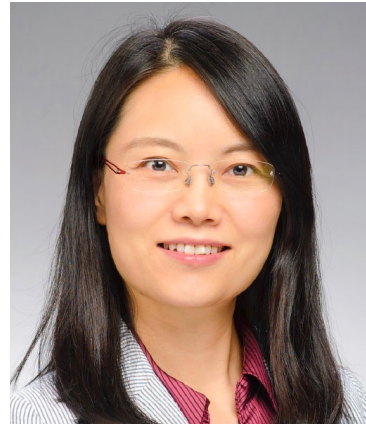
# Knowledge-enhanced Graph Learning



Yijun Tian  
ND



Shichao Pei  
UMB



Xiangliang Zhang  
ND



Wei Wang  
UCLA



Hanghang Tong  
UIUC



Nitesh Chawla  
ND





AAAI-24 Tutorial  
Feb 2024, Vancouver



Tutorial Website:



[yijuntian.com/tutorial](https://yijuntian.com/tutorial)

Yijun is looking for  
job opportunities

# Tutorial Outline



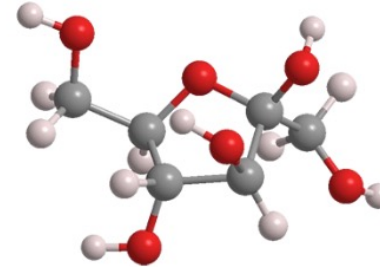
- Preliminaries and Foundations
- Graph Learning Enhanced by Knowledge from Data
- Graph Learning Enhanced by Knowledge from Models
- Graph Learning Enhanced by Knowledge from Humans and Domains
- Graph Learning Enhanced by Knowledge from External Sources
- Knowledge-enhanced Graph Learning for Real-world Applications
- Summary and Future Directions

# Graph-structured data are ubiquitous



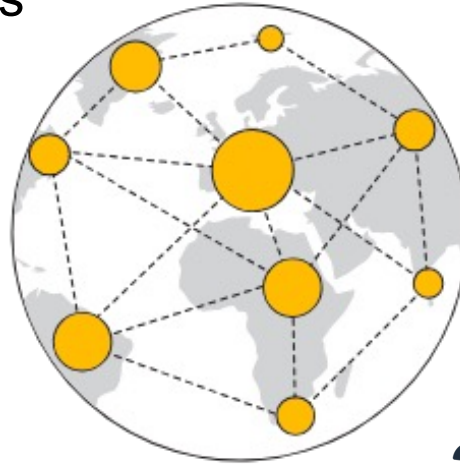
## Social Media

Nodes are people  
Edges are friendships



## Chemistry

Nodes are atoms  
Edges are chemical bonds



## Finance

Nodes are accounts  
Edges are transactions



## E-commerce

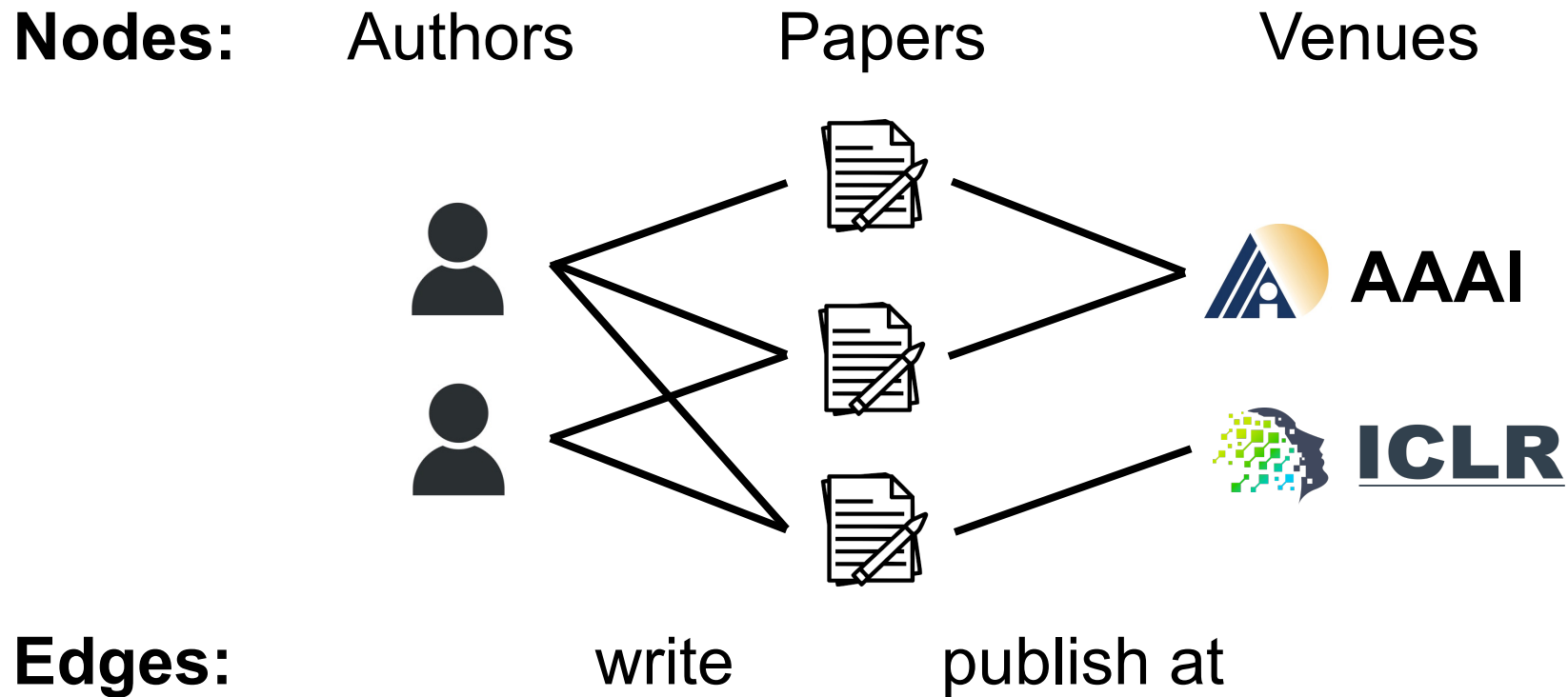
Nodes are users/products  
Edges are purchases



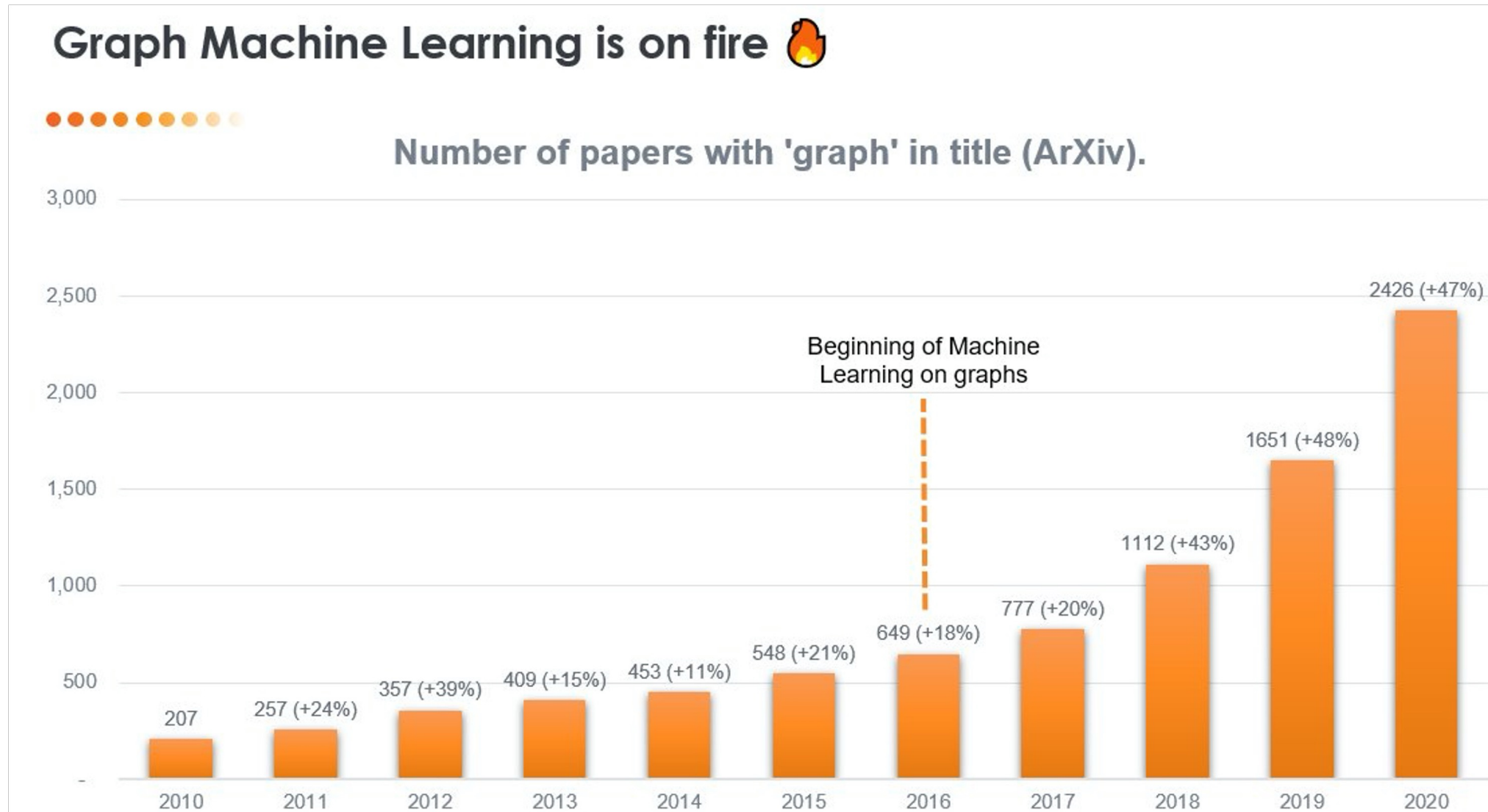
# Graph-structured data are ubiquitous



A real-world graph example:



# Graph Machine Learning: Recent Trending



# Graph Machine Learning: Previous Tutorials

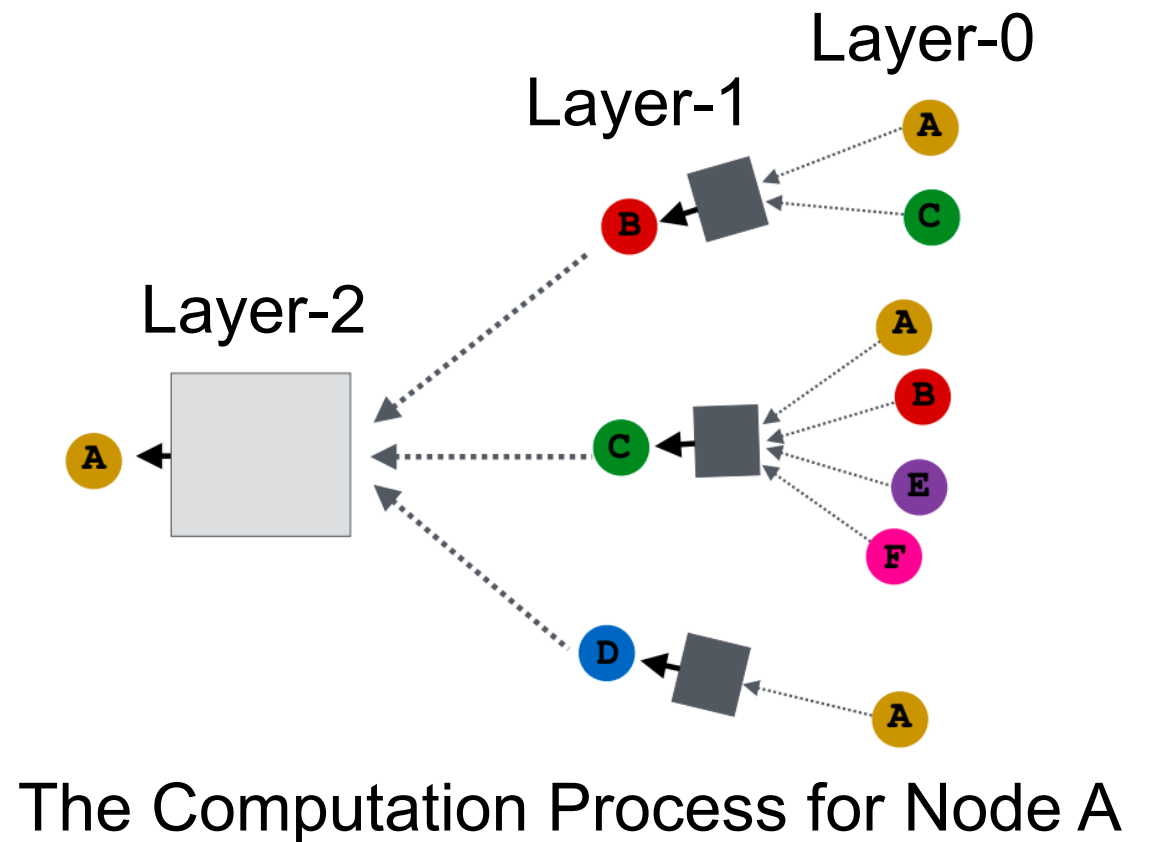
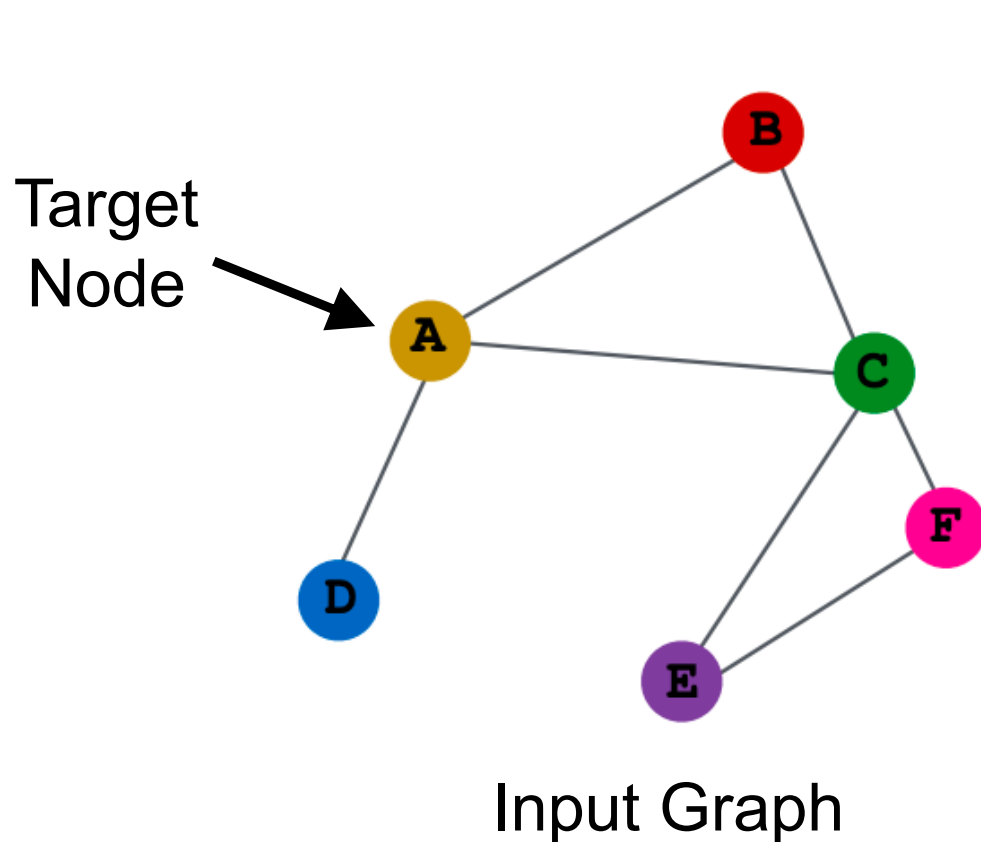


- Graph Neural Networks: Foundations, Frontiers, Applications
  - IJCAI 2022, KDD 2022, AAI 2023, KDD 2023, WWW 2023
- Graph Neural Networks: Models and Applications
  - AAI 2020, AAI 2021
- Large-Scale Graph Neural Networks: The Past and New Frontiers
  - KDD 2023
- Self-supervised Learning and Pre-training on Graphs
  - WWW 2023

# Graph Neural Networks: Foundations



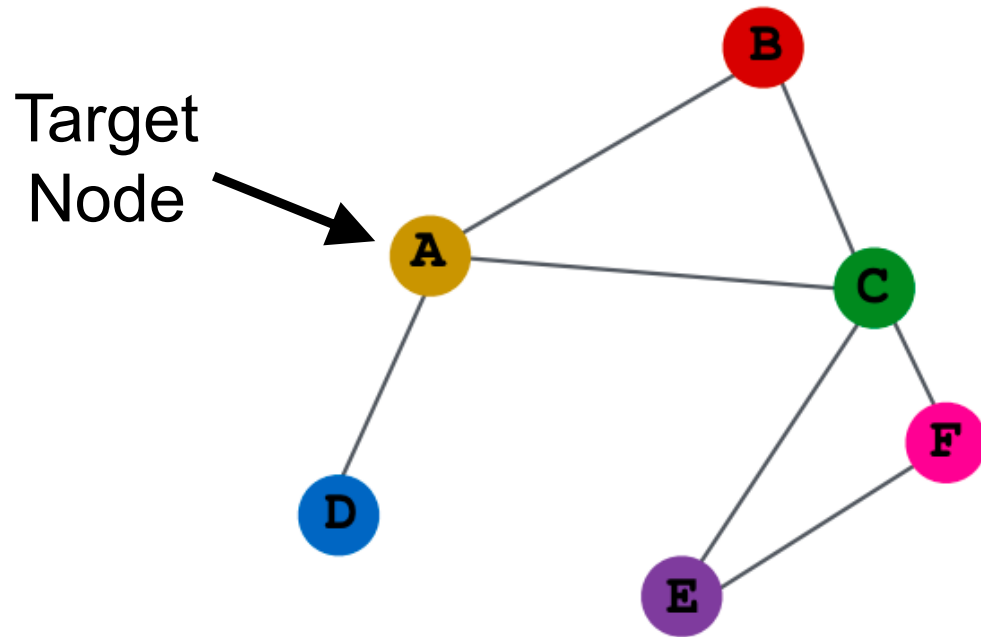
**Key idea:** aggregates information from node neighborhood



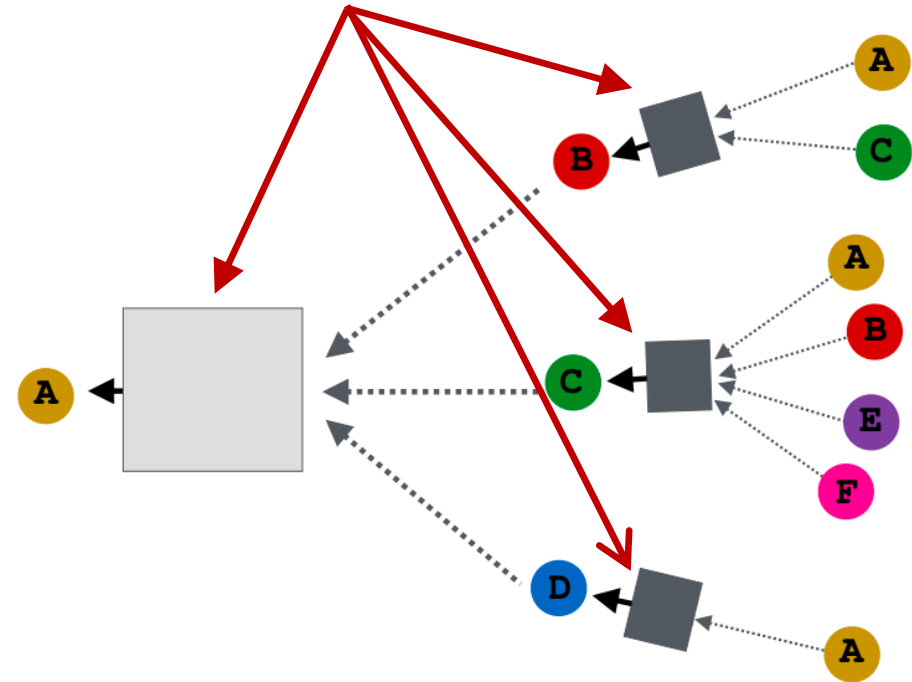
# Graph Neural Networks: Foundations



**Key idea:** aggregates information from node neighborhood



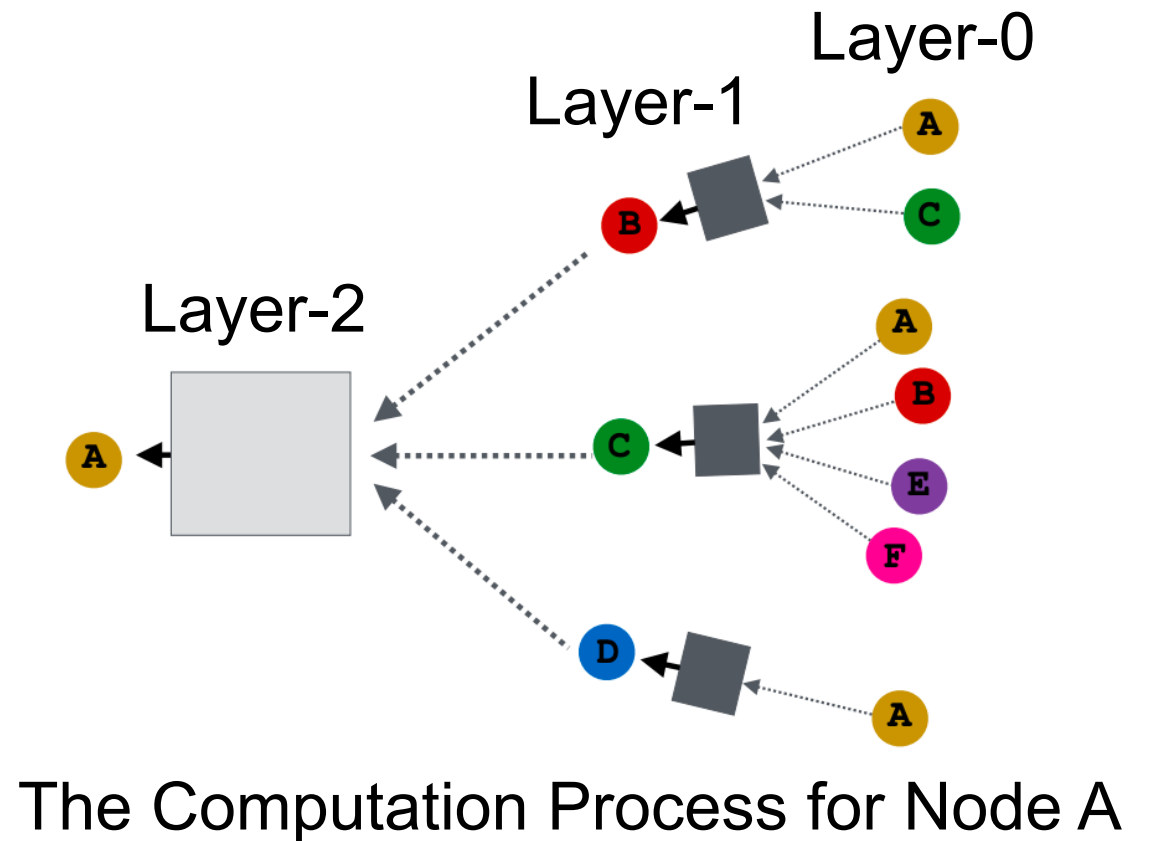
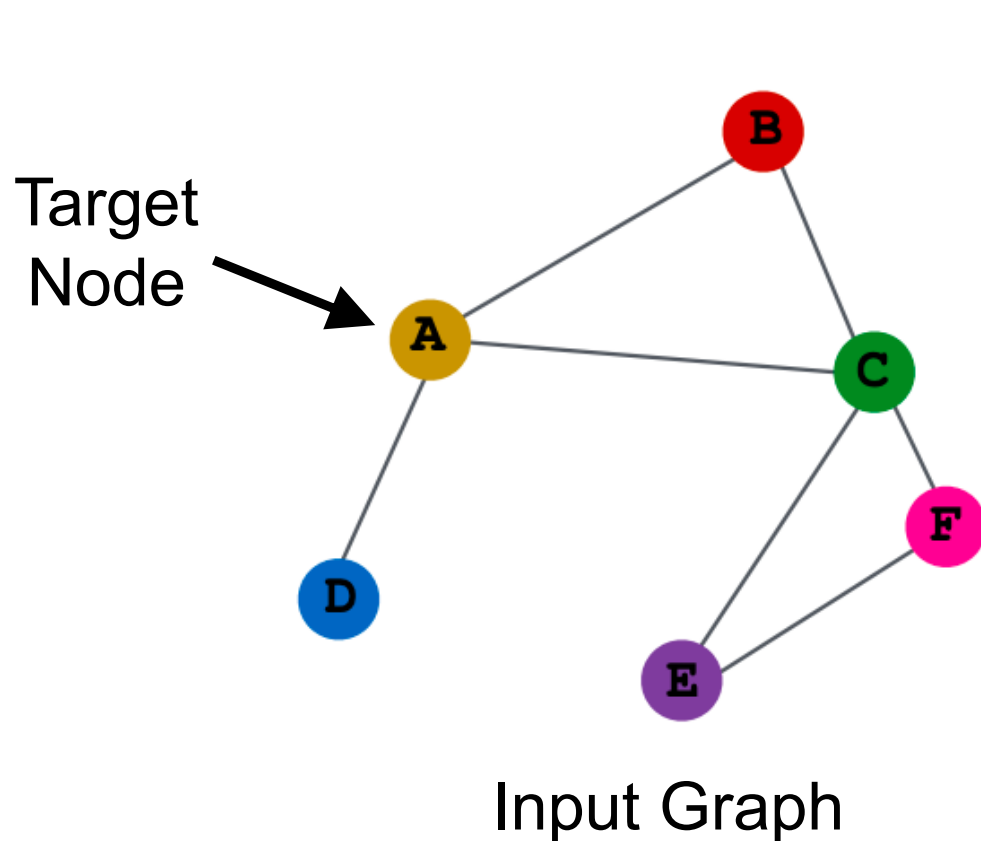
Neural networks that aggregate information



# Graph Neural Networks: Foundations



- Nodes have embeddings at each layer
- “Layer-0” embedding is the input feature of nodes

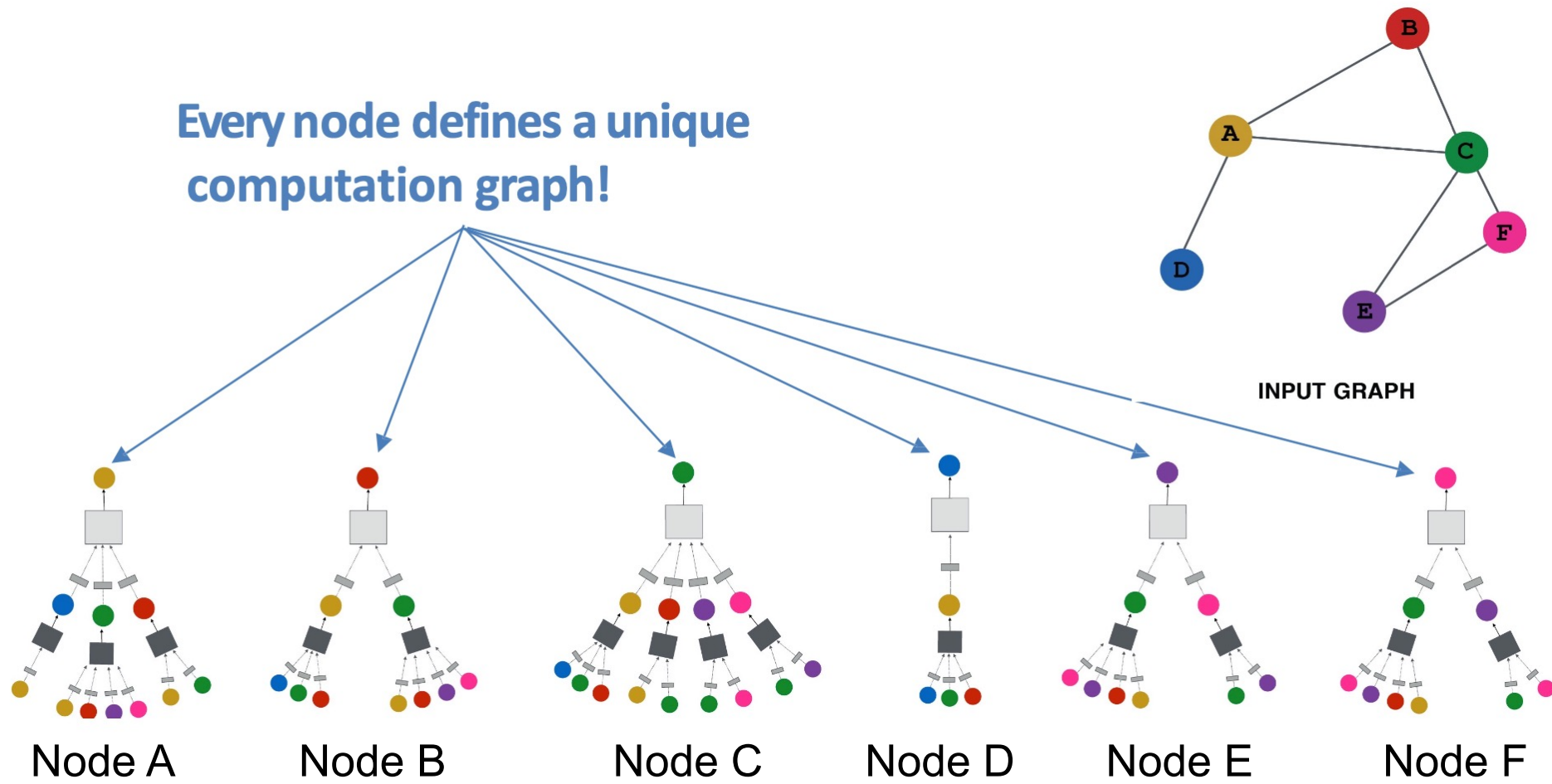




# Graph Neural Networks: Foundations



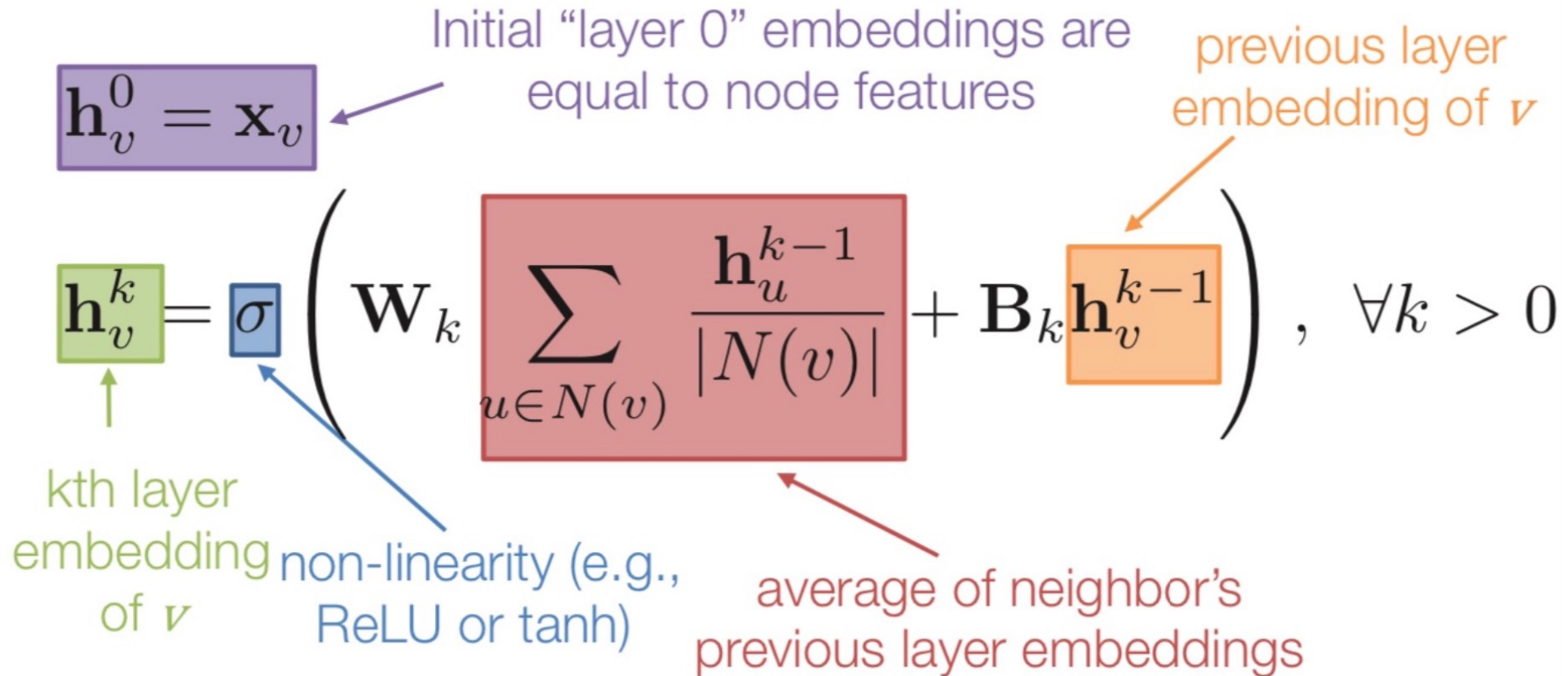
The **computation graph** is defined by neighborhood



# Graph Neural Networks: Foundations



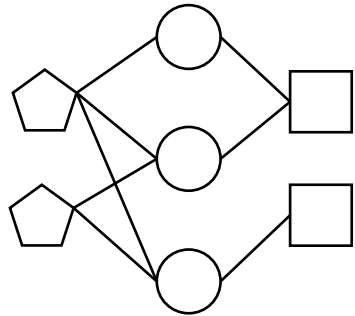
**Basic approach:** Average neighbor information and apply a neural network



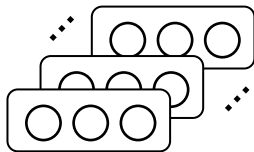
# Graph Neural Networks: Foundations



Existing GNNs primarily focus on leveraging **graph structures** and **node features**



Graph Structure:  $A$



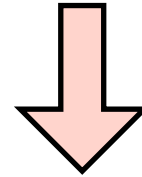
Node Features:  $X$

$$GNN = f(\text{Graph Structure}, \text{Node Features})$$

| Model     | Conference   | Citation |
|-----------|--------------|----------|
| GCN       | ICLR 2017    | 29000+   |
| GraphSAGE | NeurIPS 2017 | 12000+   |
| GAT       | ICLR 2018    | 8000+    |
| GIN       | ICLR 2019    | 6000+    |

# Graph Neural Networks: New Frontiers

Existing GNNs primarily focus on leveraging graph structures and node features



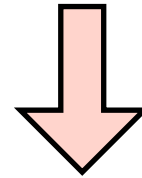
Enhancing graph learning with auxiliary knowledge



# Graph Neural Networks: New Frontiers



Existing GNNs primarily focus on leveraging graph structures and node features



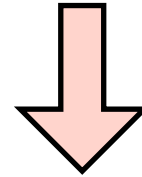
Enhancing graph learning with **auxiliary knowledge**



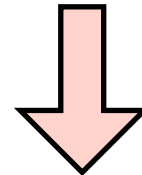
**Important and useful information** that can be obtained, extracted, or learned from resources **beyond the provided graph structures and node features**

# Graph Neural Networks: New Frontiers

Existing GNNs primarily focus on leveraging graph structures and node features



Enhancing graph learning with auxiliary knowledge



**Knowledge-enhanced Graph Learning**

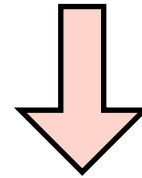




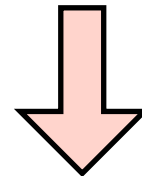
# Graph Neural Networks: New Frontiers



Existing GNNs primarily focus on leveraging graph structures and node features



Enhancing graph learning with auxiliary knowledge



**Knowledge-enhanced Graph Learning**



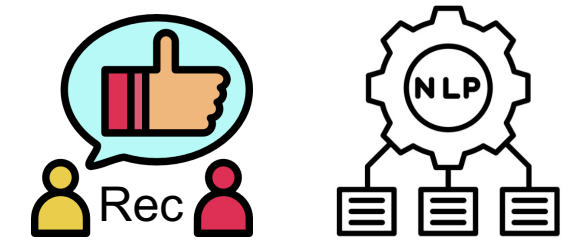
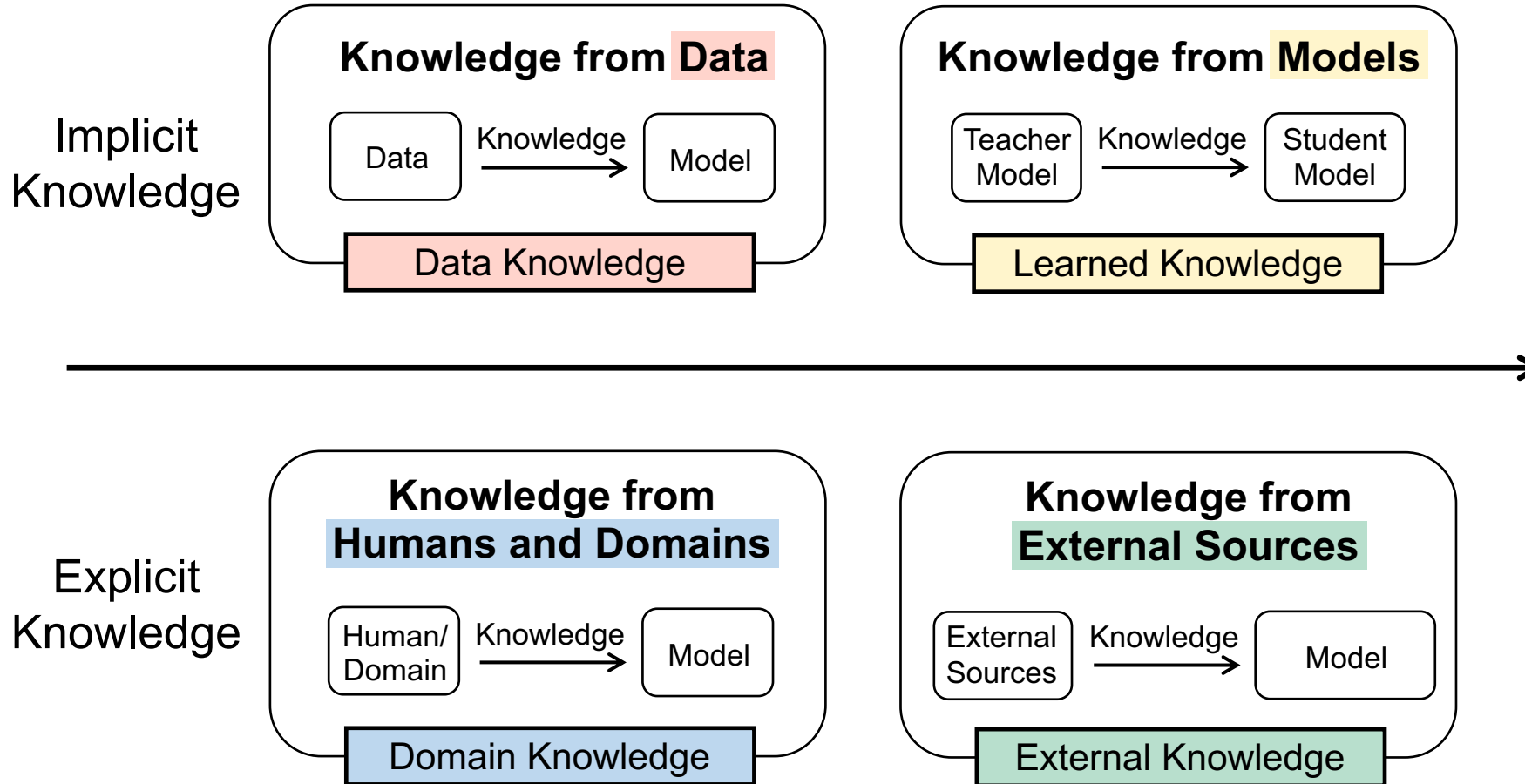
Reduce reliance on massive data and intricate model



Performance ↑  
Trustworthiness ↑  
Efficiency ↑

# Knowledge-enhanced Graph Learning

## Agenda



Real-world Applications



Future Directions

# Tutorial Outline



- Preliminaries and Foundations

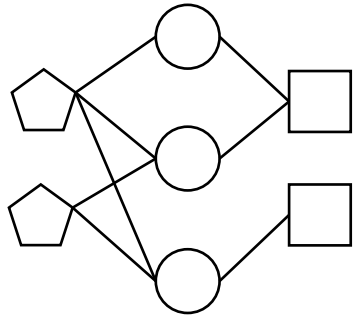


- **Graph Learning Enhanced by Knowledge from Data**
- Graph Learning Enhanced by Knowledge from Models
- Graph Learning Enhanced by Knowledge from Humans and Domains
- Graph Learning Enhanced by Knowledge from External Sources
- Knowledge-enhanced Graph Learning for Real-world Applications
- Summary and Future Directions

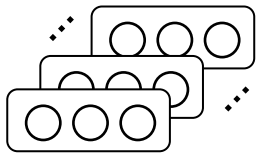
# Graph Data are Complex



## Data Property

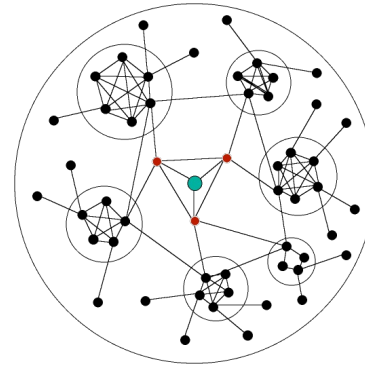


Graph Structure

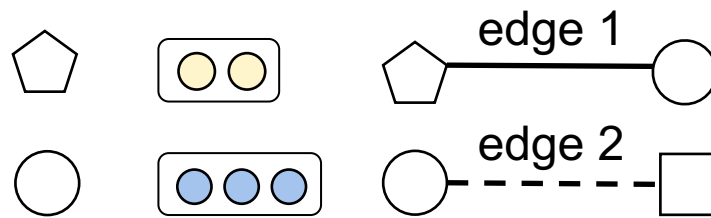


Node Features

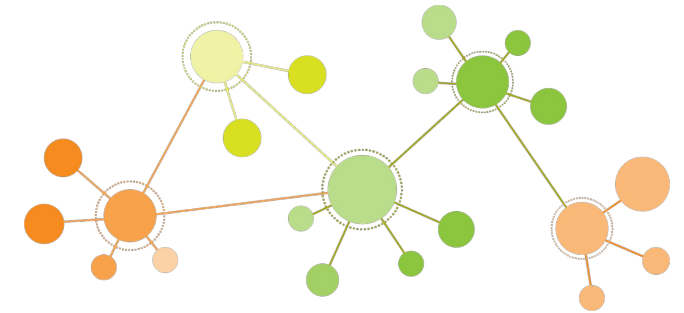
## Knowledge



Local communities



Heterogeneous Information



Node positions

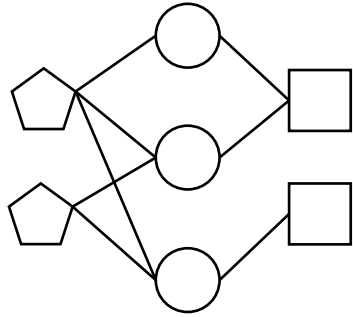
|    | p1 | p2 | p3 |           | p1 | p2 | p3 |   |
|----|----|----|----|-----------|----|----|----|---|
| p1 | 1  | 1  | 0  | Co-author | p1 | 1  | 1  | 0 |
| p2 | 1  | 1  | 1  |           | p2 | 1  | 1  | 0 |
| p3 | 0  | 1  | 1  |           | p3 | 0  | 0  | 1 |

Higher-order Semantics

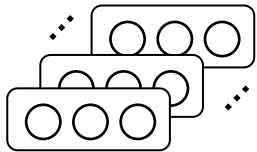
# Fundamental GNNs Focus on Data Property



## Data Property



Graph Structure



Node Features

$$GNN = f(\text{Graph Structure, Node Features})$$

| Model     | Conference   | Citation |
|-----------|--------------|----------|
| GCN       | ICLR 2017    | 29000+   |
| GraphSAGE | NeurIPS 2017 | 12000+   |
| GAT       | ICLR 2018    | 8000+    |
| GIN       | ICLR 2019    | 6000+    |

Knowledge is missed

# Taxonomy of Knowledge from Data



Knowledge can be obtained from

## 1) **single-instance** level perception

- **Node sampling** for node positions

## 2) **multiple-instance** level perception

- **Path sampling** for positional and semantic information
- **Subgraph sampling** for community information



# Taxonomy of Knowledge from Data



Knowledge can be obtained from

## 1) **single-instance** level perception



- **Node sampling for node positions**

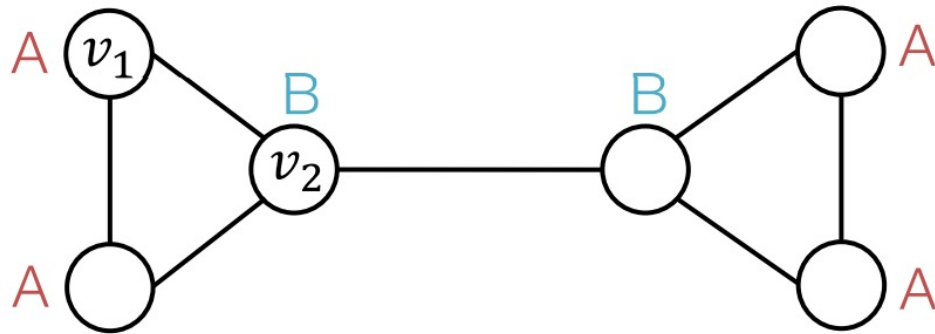
## 2) **multiple-instance** level perception

- Path sampling for positional and semantic information
- Subgraph sampling for community information

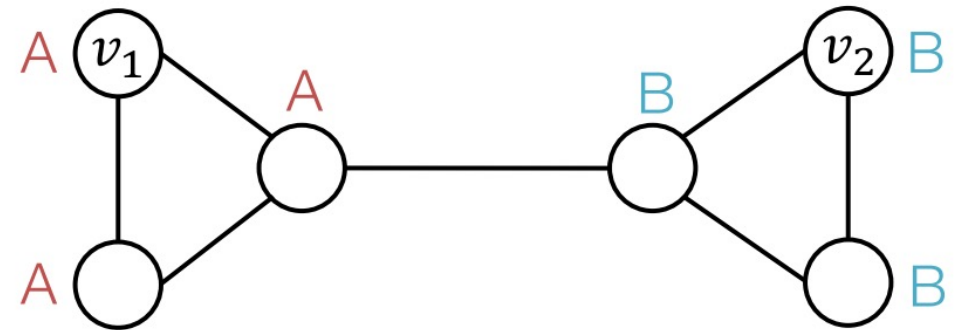
# Node Sampling for Node Positions

Why do we need to emphasize node position?

Different types of tasks



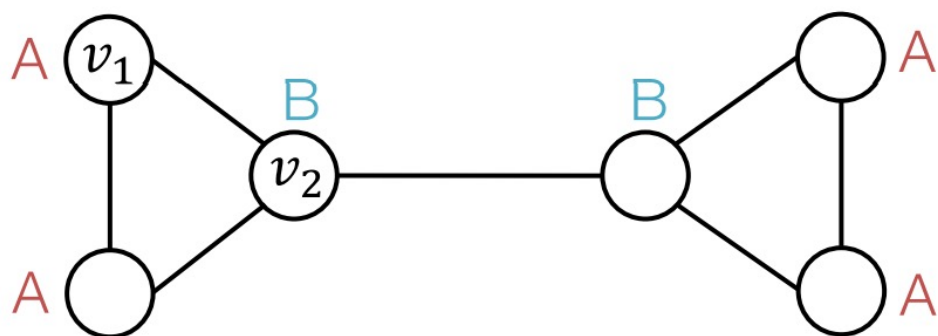
Structure-aware task: nodes labeled by **structural roles**



Position-aware task: nodes labeled by **positions**

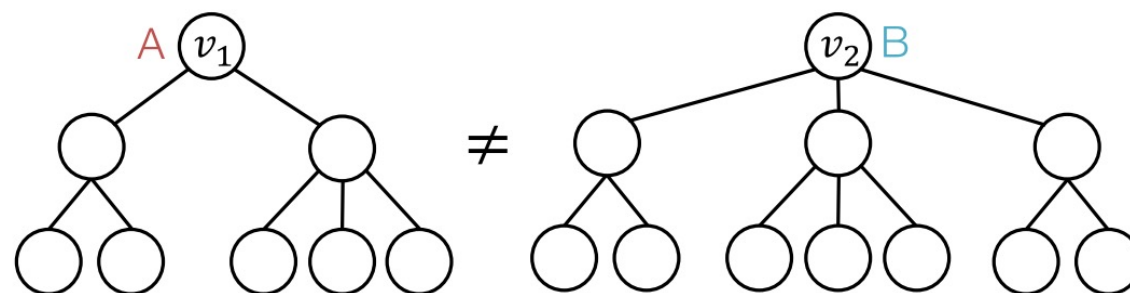
# Node Sampling for Node Positions

GNNs usually work well on **structure-aware task**



Structure-aware task: nodes labeled by **structural roles**

Computation graph of nodes  $v_1$  and  $v_2$  differs

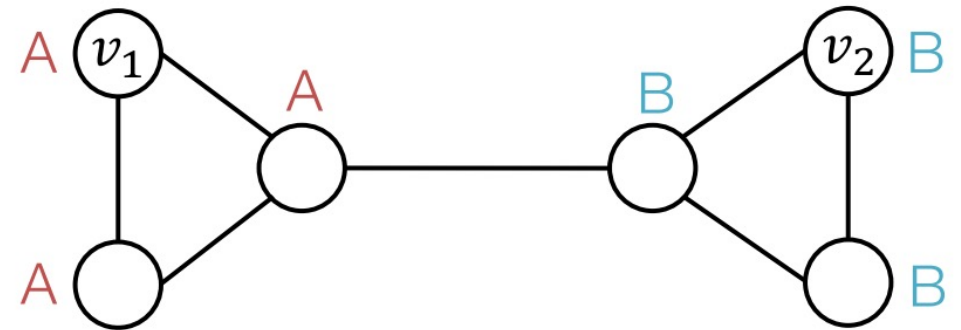
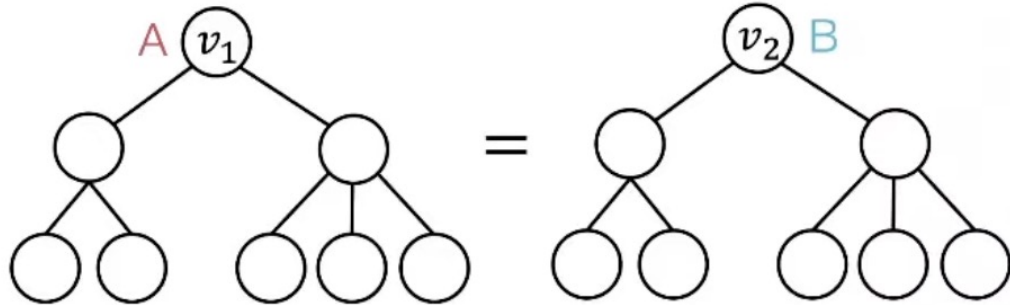


# Node Sampling for Node Positions



GNNs usually perform poorly on **position-aware task**

Computation graph of nodes  $v_1$  and  $v_2$  are the same



Position-aware task: nodes labeled by **positions**

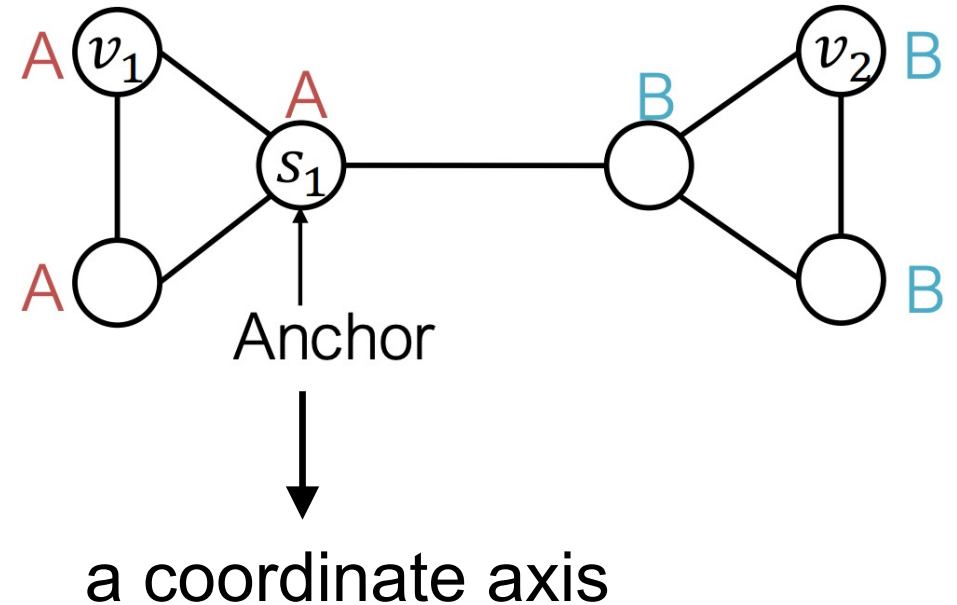
# Node Sampling for Node Positions

To encode node positions, **define anchor**

- Randomly **pick a node** as an anchor node
- Represent  $v_1$  and  $v_2$  with shortest distances to the anchor

↙

|       |       |
|-------|-------|
|       | $s_1$ |
| $v_1$ | 1     |
| $v_2$ | 2     |



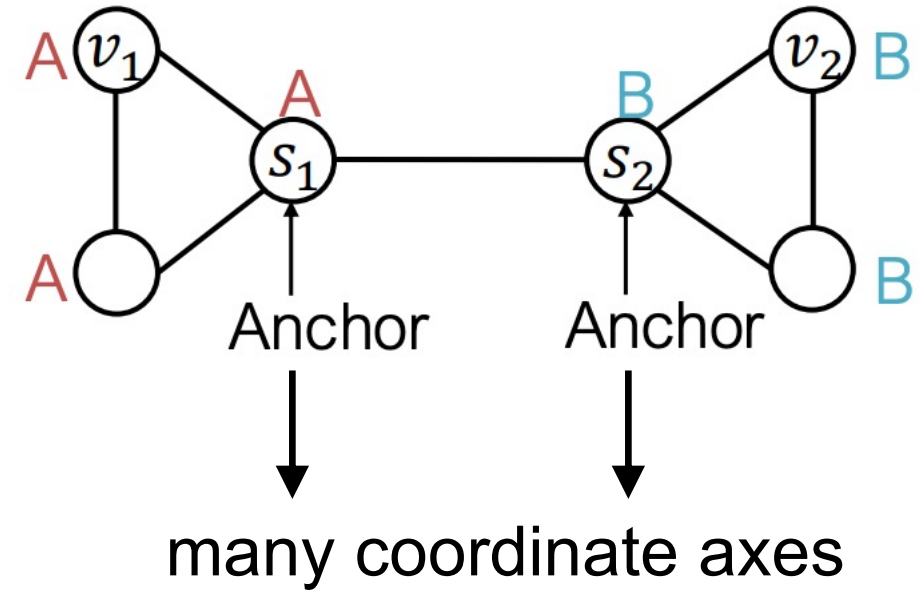
# Node Sampling for Node Positions

To encode node positions, **define anchor**

- pick **more nodes** as anchors for more coordinate axes
- Represent  $v_1$  and  $v_2$  with shortest distances to these anchors

→

|       | $s_1$ | $s_2$ |
|-------|-------|-------|
| $v_1$ | 1     | 2     |
| $v_2$ | 2     | 1     |



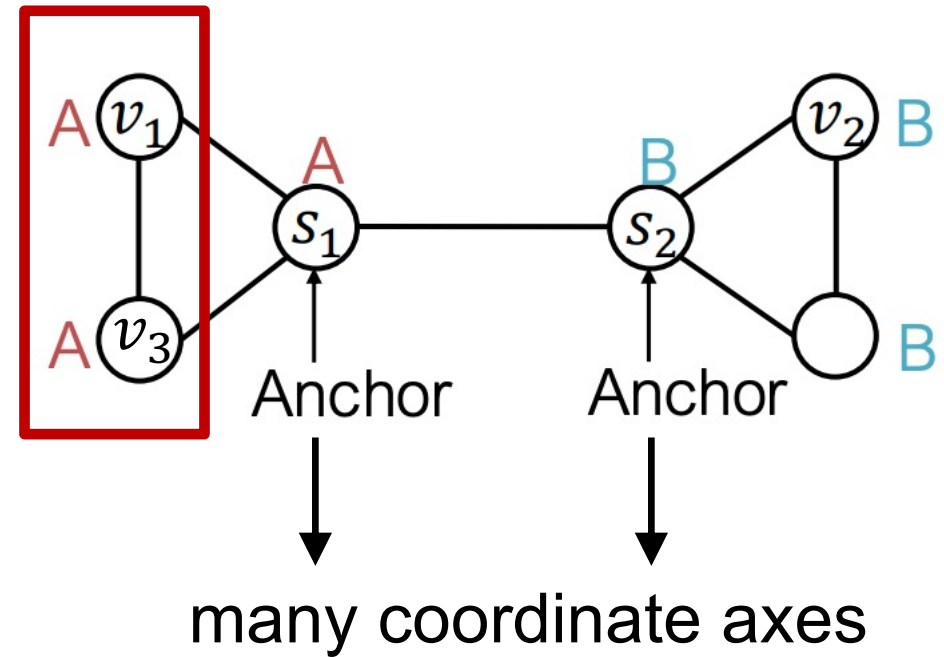


# Node Sampling for Node Positions

To encode node positions, **define anchor**

- Problem:  $v_1$  and  $v_3$  have same shortest distances to these anchors

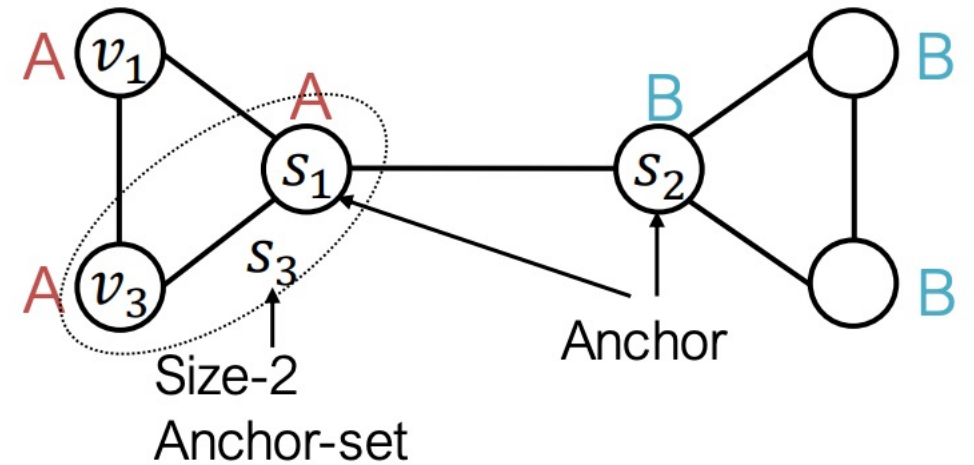
|       | $s_1$ | $s_2$ |
|-------|-------|-------|
| $v_1$ | 1     | 2     |
| $v_3$ | 1     | 2     |



# Node Sampling for Node Positions

To encode node positions, **define anchor**

- Include **multiple nodes in an anchor**: anchor-set
- Better position estimation/encoding



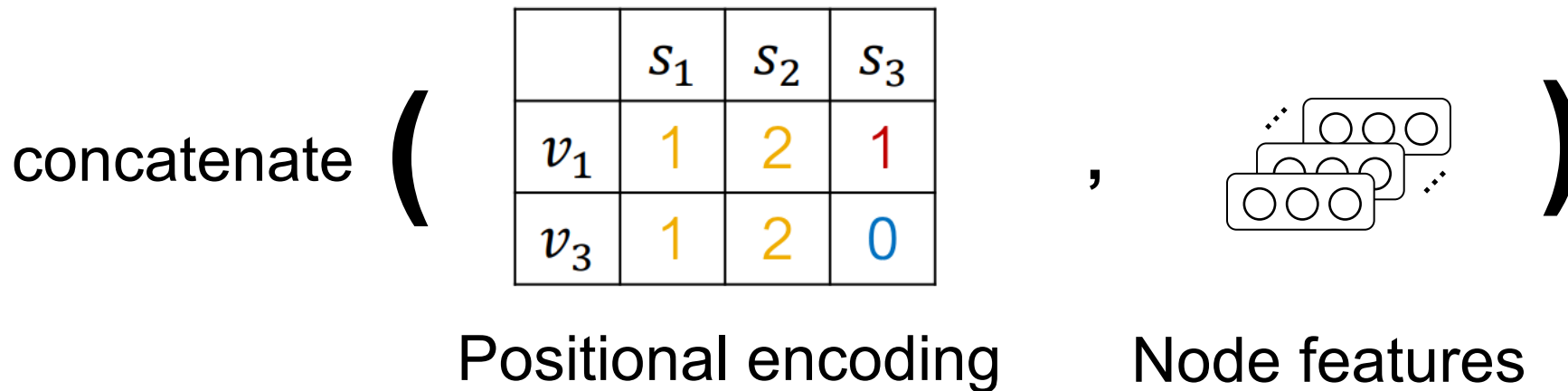
|       | $s_1$ | $s_2$ | $s_3$ |
|-------|-------|-------|-------|
| $v_1$ | 1     | 2     | 1     |
| $v_3$ | 1     | 2     | 0     |

Anchor-set can distinguish  $v_1$  and  $v_3$ , while anchor cannot

# Node Sampling for Node Positions

To use obtained positional encoding

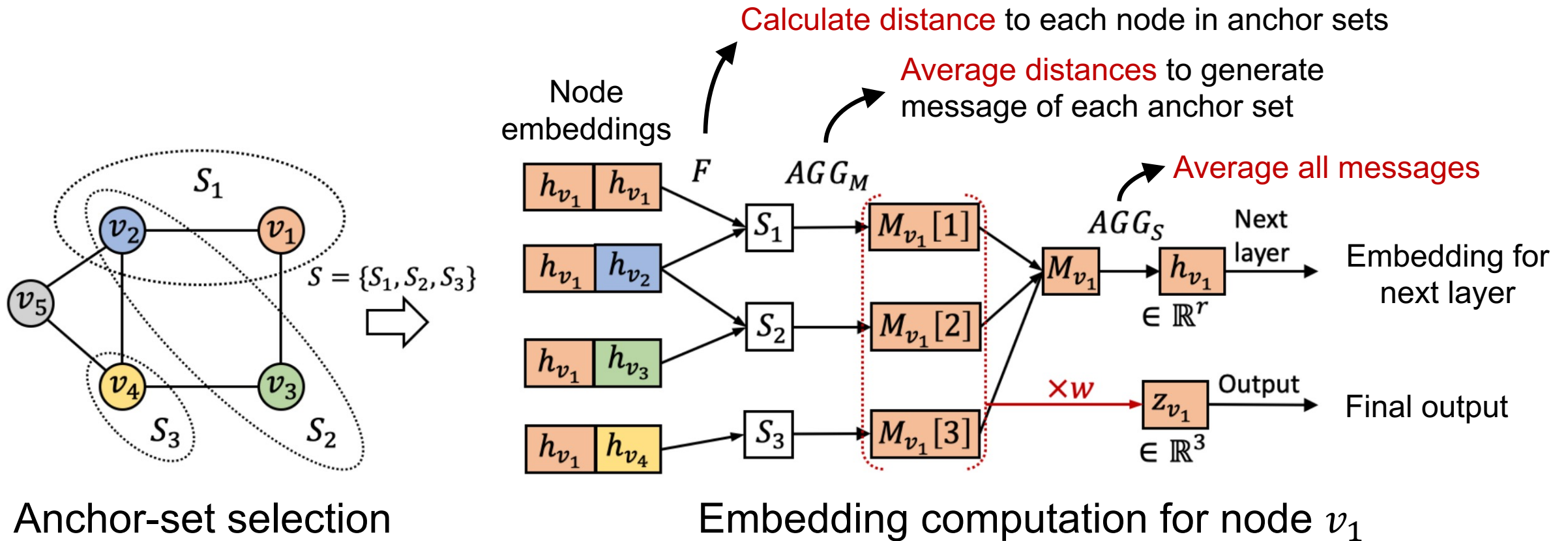
We can simply concatenate it with node features and use them as usual



Problem: sending this concatenated feature to a neural network cannot preserve the **permutation invariant property** of positional encoding

# Node Sampling for Node Positions

To use obtained positional encoding, design **P-GNN** with permutation invariant aggregation (e.g., mean)



# Node Sampling for Node Positions

P-GNNs outperforms GNNs



Link Prediction task, measured in ROC AUC

|            | Grid-T               | Communities-T        | Grid                 | Communities          | PPI                  |
|------------|----------------------|----------------------|----------------------|----------------------|----------------------|
| GCN        | 0.698 ± 0.051        | 0.981 ± 0.004        | 0.456 ± 0.037        | 0.512 ± 0.008        | 0.769 ± 0.002        |
| GraphSAGE  | 0.682 ± 0.050        | 0.978 ± 0.003        | 0.532 ± 0.050        | 0.516 ± 0.010        | 0.803 ± 0.005        |
| GAT        | 0.704 ± 0.050        | 0.980 ± 0.005        | 0.566 ± 0.052        | 0.618 ± 0.025        | 0.783 ± 0.004        |
| GIN        | 0.732 ± 0.050        | 0.984 ± 0.005        | 0.499 ± 0.054        | 0.692 ± 0.049        | 0.782 ± 0.010        |
| P-GNN-F-1L | 0.542 ± 0.057        | 0.930 ± 0.093        | 0.619 ± 0.080        | 0.939 ± 0.083        | 0.719 ± 0.027        |
| P-GNN-F-2L | 0.637 ± 0.078        | <b>0.989</b> ± 0.003 | 0.694 ± 0.066        | <b>0.991</b> ± 0.003 | 0.805 ± 0.003        |
| P-GNN-E-1L | 0.665 ± 0.033        | 0.966 ± 0.013        | 0.879 ± 0.039        | 0.985 ± 0.005        | 0.775 ± 0.029        |
| P-GNN-E-2L | <b>0.834</b> ± 0.099 | 0.988 ± 0.003        | <b>0.940</b> ± 0.027 | 0.985 ± 0.008        | <b>0.808</b> ± 0.003 |

Different variants of P-GNN with fast or regular way of calculating shortest path

+66%

When graphs come with rich features (e.g., PPI dataset), the performance improvement is smaller, because node features may already capture positional information

# Taxonomy of Knowledge from Data



Knowledge can be obtained from

## 1) **single-instance** level perception

- Node sampling for node positions

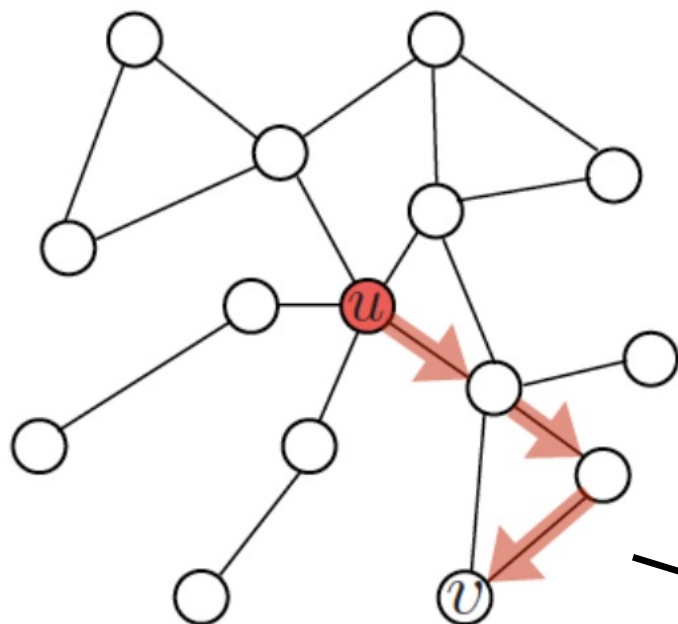
## 2) **multiple-instance** level perception



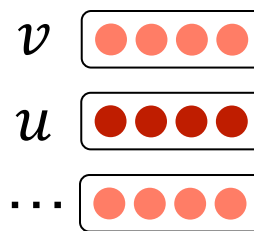
- **Path sampling for positional and semantic information**
- Subgraph sampling for community information

# Path Sampling for Positions and Semantics

Sample a path of nodes for **positions**, e.g., via random walk



- Nodes appeared in a sampled path are similar
- Similarity between nodes in a graph translated to closeness in embedding

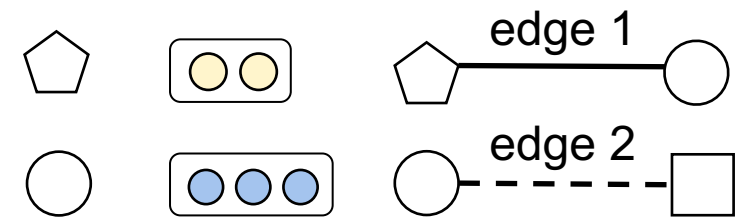
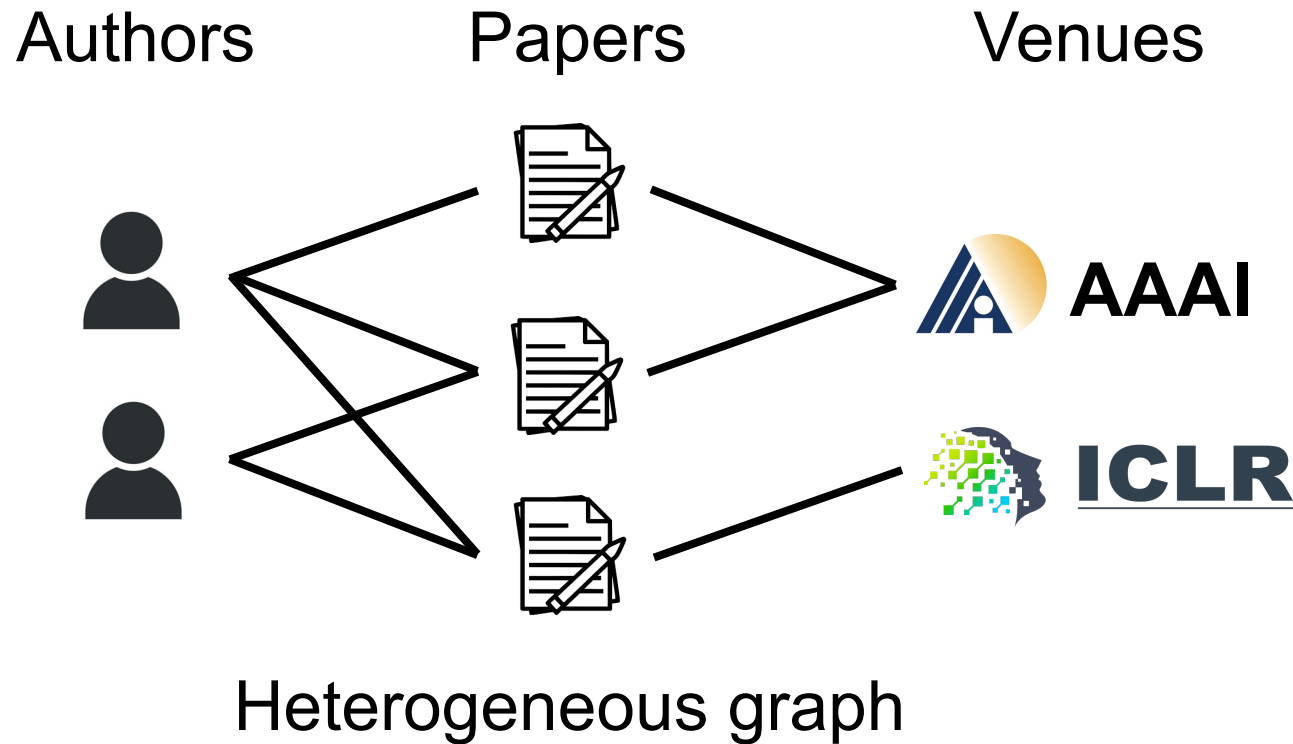


Positional embedding

# Path Sampling for Positions and Semantics



Heterogeneous graphs contain rich information with **semantics**



Heterogeneous Information



# Path Sampling for Positions and Semantics

Sample a path of nodes for **semantics**, e.g., via metapaths



Authors

Papers

Venues



AAAI



ICLR



a1



paper



a2

Co-author



p1



venue



p2

Co-venue

Heterogeneous graph

metapaths

# Path Sampling for Positions and Semantics

Sample a path of nodes for **semantics**, e.g., via metapaths



# Path Sampling for Positions and Semantics

To encode the positions and semantics, introducing **HGMAE**



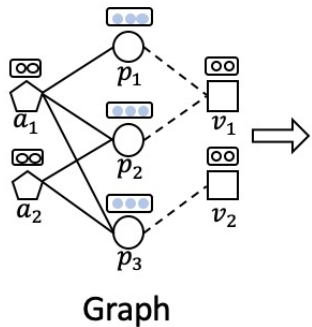
Data

Knowledge  
Extraction

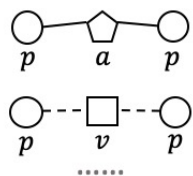
Learning Strategies

Regular GNN

Learning Objectives



Graph

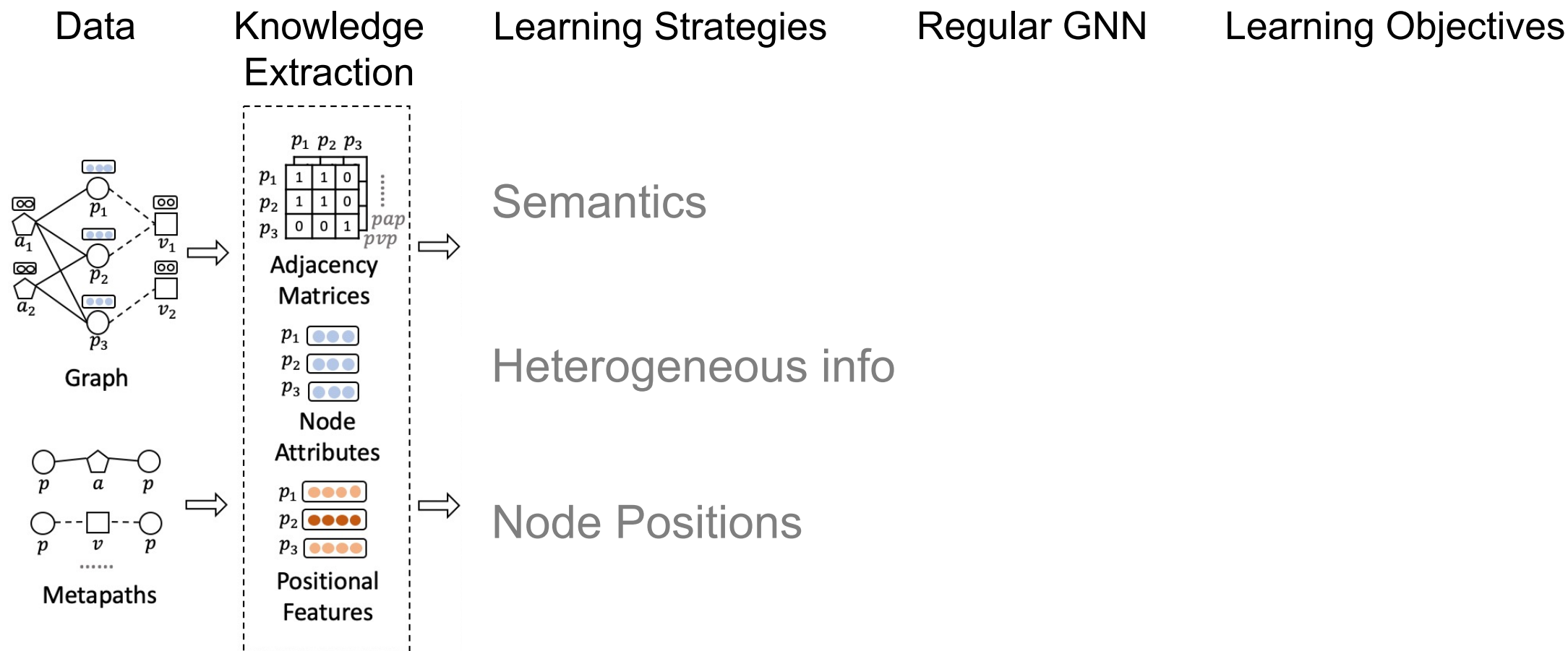


Metapaths

Metapaths involve the semantic knowledge

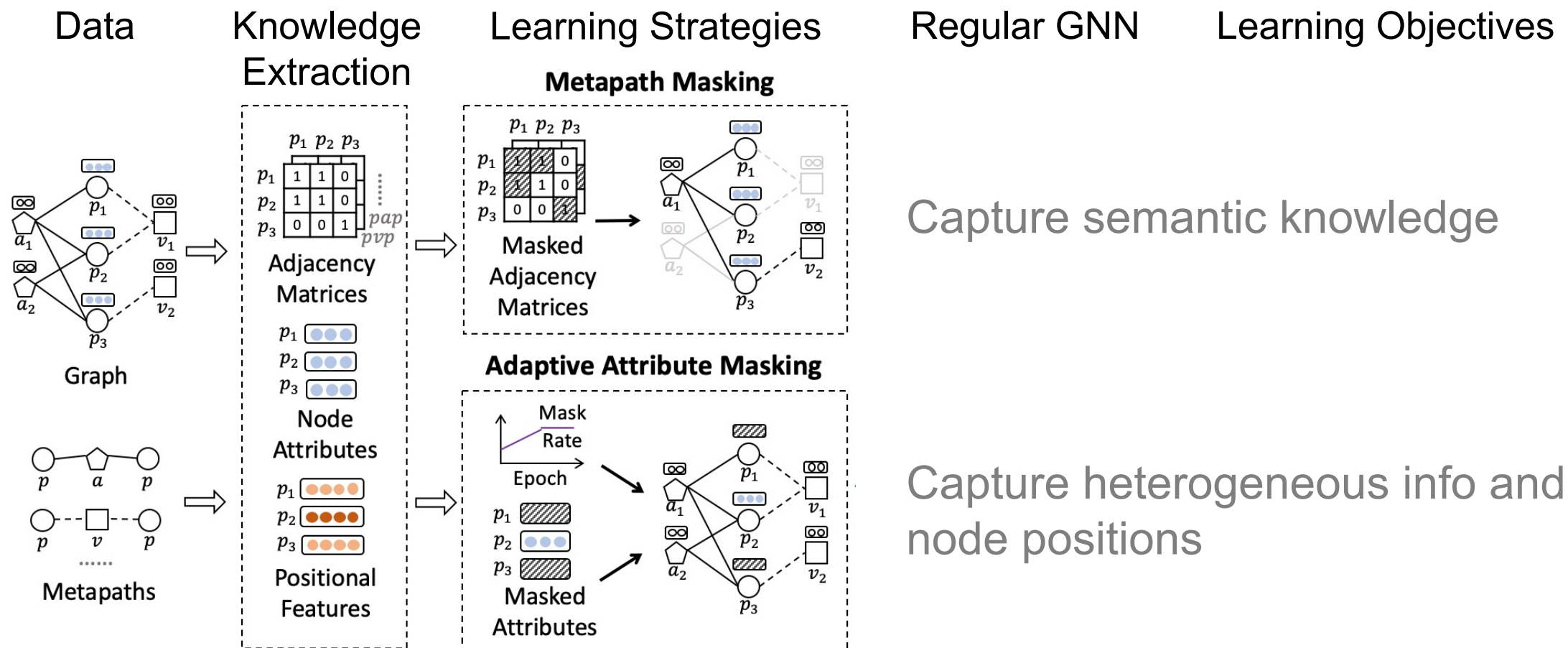
# Path Sampling for Positions and Semantics

To encode the positions and semantics, introducing **HGMAE**



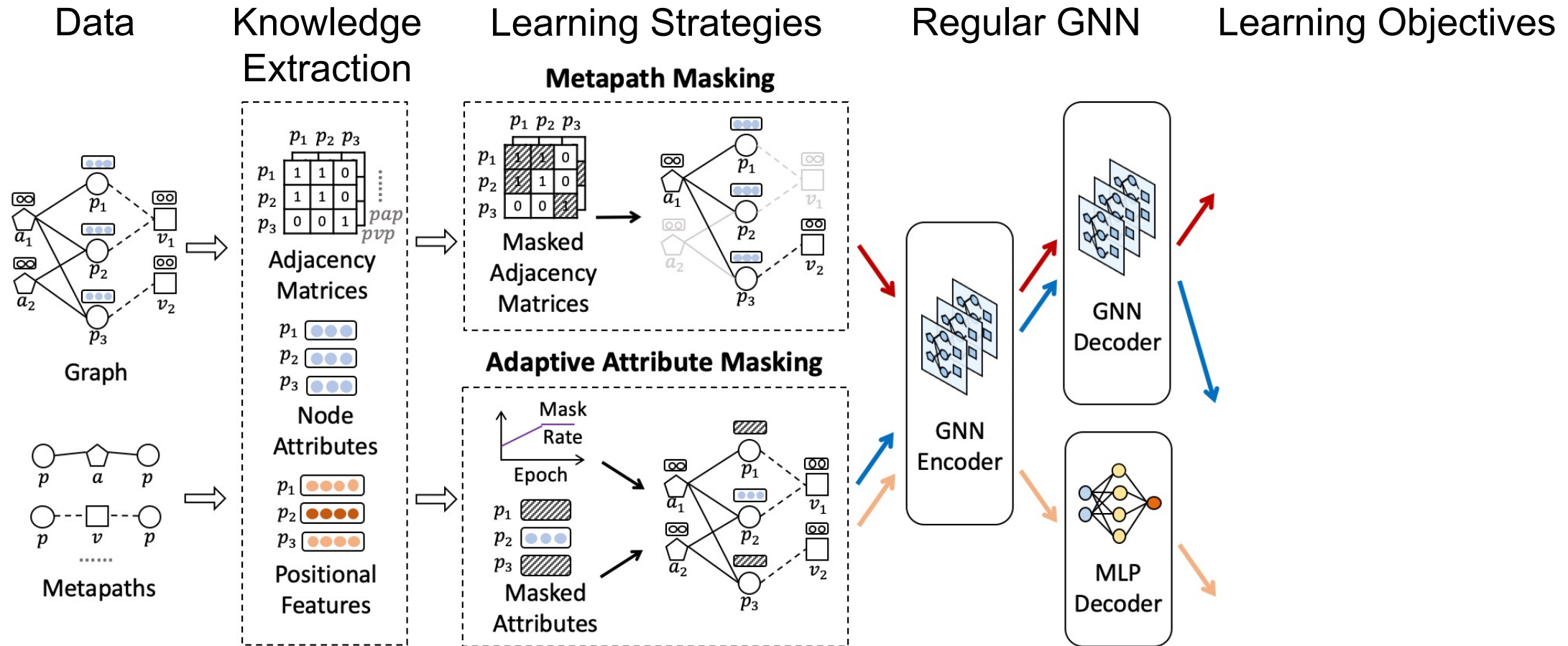
# Path Sampling for Positions and Semantics

To encode the positions and semantics, introducing **HGMAE**



# Path Sampling for Positions and Semantics

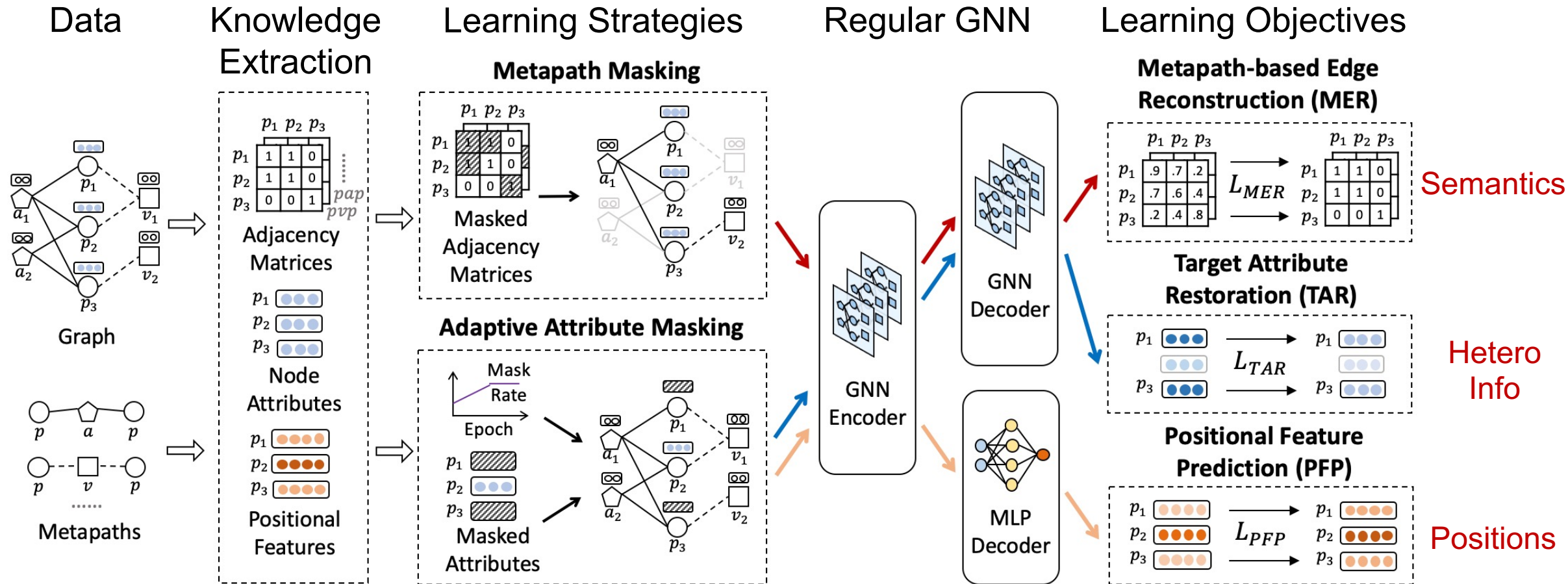
To encode the positions and semantics, introducing **HGMAE**





# Path Sampling for Positions and Semantics

To encode the positions and semantics, introducing **HGMAE**



# Path Sampling for Positions and Semantics

To encode the positions and semantics, introducing **HGMAE**



## Performance Comparison: Node Classification

| Datasets | Metric | Split | GraphSAGE | GAE       | Mp2vec    | HERec     | HetGNN    | HAN       | DGI       | DMGI      | HeCo      | GraphMAE  | HGMAE            |
|----------|--------|-------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|------------------|
| DBLP     | Mi-F1  | 20    | 71.44±8.7 | 91.55±0.1 | 89.67±0.1 | 90.24±0.4 | 90.11±1.0 | 90.16±0.9 | 88.72±2.6 | 90.78±0.3 | 91.97±0.2 | 89.31±0.7 | <b>92.71±0.5</b> |
|          |        | 40    | 73.61±8.6 | 90.00±0.3 | 89.14±0.2 | 90.15±0.4 | 89.03±0.7 | 89.47±0.9 | 89.22±0.5 | 89.92±0.4 | 90.76±0.3 | 87.80±0.5 | <b>92.43±0.3</b> |
|          |        | 60    | 74.05±8.3 | 90.95±0.2 | 91.17±0.1 | 91.01±0.3 | 90.43±0.6 | 90.34±0.8 | 90.35±0.8 | 90.66±0.5 | 91.59±0.2 | 89.82±0.4 | <b>93.05±0.3</b> |
|          | Ma-F1  | 20    | 71.97±8.4 | 90.90±0.1 | 88.98±0.2 | 89.57±0.4 | 89.51±1.1 | 89.31±0.9 | 87.93±2.4 | 89.94±0.4 | 91.28±0.2 | 87.94±0.7 | <b>92.28±0.5</b> |
|          |        | 40    | 73.69±8.4 | 89.60±0.3 | 88.68±0.2 | 89.73±0.4 | 88.61±0.8 | 88.87±1.0 | 88.62±0.6 | 89.25±0.4 | 90.34±0.3 | 86.85±0.7 | <b>92.12±0.3</b> |
|          |        | 60    | 73.86±8.1 | 90.08±0.2 | 90.25±0.1 | 90.18±0.3 | 89.56±0.5 | 89.20±0.8 | 89.19±0.9 | 89.46±0.6 | 90.64±0.3 | 88.07±0.6 | <b>92.33±0.3</b> |
|          | AUC    | 20    | 90.59±4.3 | 98.15±0.1 | 97.69±0.0 | 98.21±0.2 | 97.96±0.4 | 98.07±0.6 | 96.99±1.4 | 97.75±0.3 | 98.32±0.1 | 92.23±3.0 | <b>98.90±0.1</b> |
|          |        | 40    | 91.42±4.0 | 97.85±0.1 | 97.08±0.0 | 97.93±0.1 | 97.70±0.3 | 97.48±0.6 | 97.12±0.4 | 97.23±0.2 | 98.06±0.1 | 91.76±2.5 | <b>98.55±0.1</b> |
|          |        | 60    | 91.73±3.8 | 98.37±0.1 | 98.00±0.0 | 98.49±0.1 | 97.97±0.2 | 97.96±0.5 | 97.76±0.5 | 97.72±0.4 | 98.59±0.1 | 91.63±2.5 | <b>98.89±0.1</b> |
| Freebase | Mi-F1  | 20    | 54.83±3.0 | 55.20±0.7 | 56.23±0.8 | 57.92±0.5 | 56.85±0.9 | 57.24±3.2 | 58.16±0.9 | 58.26±0.9 | 61.72±0.6 | 64.88±1.8 | <b>65.15±1.3</b> |
|          |        | 40    | 57.08±3.2 | 56.05±2.0 | 61.01±1.3 | 62.71±0.7 | 53.96±1.1 | 63.74±2.7 | 57.82±0.8 | 54.28±1.6 | 64.03±0.7 | 62.34±1.0 | <b>67.23±0.8</b> |
|          |        | 60    | 55.92±3.2 | 53.85±0.4 | 58.74±0.8 | 58.57±0.5 | 56.84±0.7 | 61.06±2.0 | 57.96±0.7 | 56.69±1.2 | 63.61±1.6 | 59.48±6.2 | <b>67.44±1.2</b> |
|          | Ma-F1  | 20    | 45.14±4.5 | 53.81±0.6 | 53.96±0.7 | 55.78±0.5 | 52.72±1.0 | 53.16±2.8 | 54.90±0.7 | 55.79±0.9 | 59.23±0.7 | 59.04±1.0 | <b>62.06±1.0</b> |
|          |        | 40    | 44.88±4.1 | 52.44±2.3 | 57.80±1.1 | 59.28±0.6 | 48.57±0.5 | 59.63±2.3 | 53.40±1.4 | 49.88±1.9 | 61.19±0.6 | 56.40±1.1 | <b>64.64±0.9</b> |
|          |        | 60    | 45.16±3.1 | 50.65±0.4 | 55.94±0.7 | 56.50±0.4 | 52.37±0.8 | 56.77±1.7 | 53.81±1.1 | 52.10±0.7 | 60.13±1.3 | 51.73±2.3 | <b>63.84±1.0</b> |
|          | AUC    | 20    | 67.63±5.0 | 73.03±0.7 | 71.78±0.7 | 73.89±0.4 | 70.84±0.7 | 73.26±2.1 | 72.80±0.6 | 73.19±1.2 | 76.22±0.8 | 72.60±0.2 | <b>78.36±1.1</b> |
|          |        | 40    | 66.42±4.7 | 74.05±0.9 | 75.51±0.8 | 76.08±0.4 | 69.48±0.2 | 77.74±1.2 | 72.97±1.1 | 70.77±1.6 | 78.44±0.5 | 72.44±1.6 | <b>79.69±0.7</b> |
|          |        | 60    | 66.78±3.5 | 71.75±0.4 | 74.78±0.4 | 74.89±0.4 | 71.01±0.5 | 75.69±1.5 | 73.32±0.9 | 73.17±1.4 | 78.04±0.4 | 70.66±1.6 | <b>79.11±1.3</b> |

+7.5%



# Path Sampling for Positions and Semantics

To encode the positions and semantics, introducing **HGMAE**



## Performance Comparison: Node Clustering

| Datasets     | DBLP         |              | Freebase     |              | ACM          |              | AMiner       |              |
|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Metrics      | NMI          | ARI          | NMI          | ARI          | NMI          | ARI          | NMI          | ARI          |
| GraphSage    | 51.50        | 36.40        | 9.05         | 10.49        | 29.20        | 27.72        | 15.74        | 10.10        |
| GAE          | 72.59        | 77.31        | 19.03        | 14.10        | 27.42        | 24.49        | 28.58        | 20.90        |
| Mp2vec       | 73.55        | 77.70        | 16.47        | 17.32        | 48.43        | 34.65        | 30.80        | 25.26        |
| HERec        | 70.21        | 73.99        | 19.76        | 19.36        | 47.54        | 35.67        | 27.82        | 20.16        |
| HetGNN       | 69.79        | 75.34        | 12.25        | 15.01        | 41.53        | 34.81        | 21.46        | 26.60        |
| DGI          | 59.23        | 61.85        | 18.34        | 11.29        | 51.73        | 41.16        | 22.06        | 15.93        |
| DMGI         | 70.06        | 75.46        | 16.98        | 16.91        | 51.66        | 46.64        | 19.24        | 20.09        |
| GraphMAE     | 65.86        | 69.75        | 19.43        | 20.05        | 47.03        | 46.48        | 17.98        | 21.52        |
| HeCo         | 74.51        | 80.17        | 20.38        | 20.98        | 56.87        | 56.94        | 32.26        | 28.64        |
| <b>HGMAE</b> | <b>76.92</b> | <b>82.34</b> | <b>22.05</b> | <b>22.84</b> | <b>66.68</b> | <b>71.51</b> | <b>41.10</b> | <b>38.27</b> |

+17% +25% +27% +33%

# Path Sampling for Positions and Semantics

To encode the positions and semantics, introducing **HGMAE**



## Ablation Study

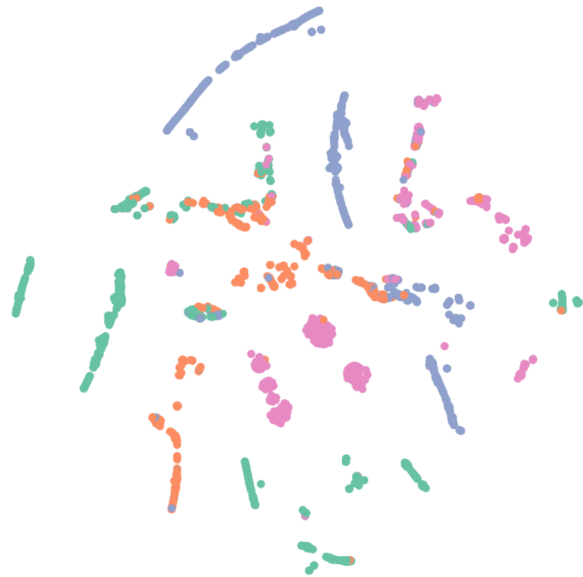
| Datasets | Metric | Semantics | Hetero Info | Positions | HGMAE            |
|----------|--------|-----------|-------------|-----------|------------------|
|          |        | w/o MER   | w/o TAR     | w/o PFP   |                  |
| ACM      | Mi-F1  | 76.85±0.2 | 88.54±0.4   | 89.81±0.5 | <b>90.59±0.5</b> |
|          | Ma-F1  | 71.93±0.4 | 88.82±0.4   | 89.94±0.4 | <b>90.80±0.5</b> |
|          | AUC    | 84.84±1.5 | 96.47±0.1   | 97.22±0.1 | <b>97.69±0.1</b> |

- Removing the learning of each knowledge results in decreased performance
- By learning all knowledge, HGMAE achieves the best results

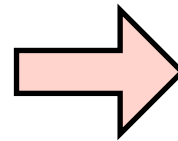
# Path Sampling for Positions and Semantics

To encode the positions and semantics, introducing **HGMAE**

How valuable is knowledge?



w/o Knowledge  
( GraphMAE )



w/ Knowledge  
( HGMAE )

# Taxonomy of Knowledge from Data



Knowledge can be obtained from

## 1) **single-instance** level perception

- Node sampling for node positions

## 2) **multiple-instance** level perception

- Path sampling for positional and semantic information
- **Subgraph sampling for community information**



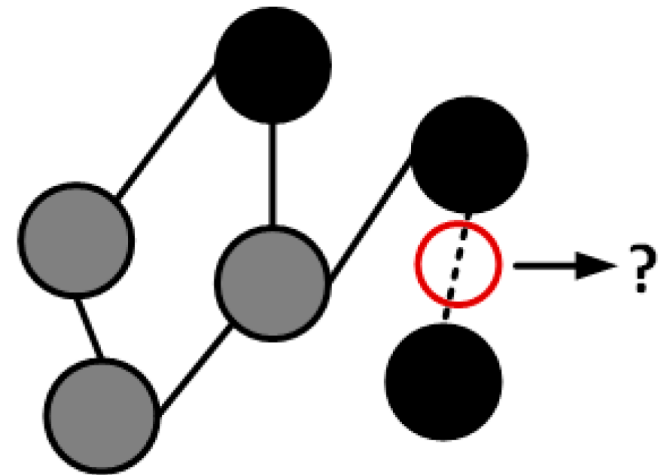
# Subgraph Sampling for Community Information



What is link prediction?

**Predicting missing or future links** between pairs of nodes is important in social networks, recommender systems, biological networks, and academic graphs

- Friend Recommendation
- Product Recommendation
- Protein-protein Interaction
- Co-authorship Prediction
- ...



# Subgraph Sampling for Community Information

## Heuristic methods for link prediction



| Name                    | Formula   | Order  |
|-------------------------|---|--------|
| common neighbors        | $ \Gamma(x) \cap \Gamma(y) $  | first  |
| Jaccard                 | $\frac{ \Gamma(x) \cap \Gamma(y) }{ \Gamma(x) \cup \Gamma(y) }$   | first  |
| preferential attachment | $ \Gamma(x)  \cdot  \Gamma(y) $   | first  |
| Adamic-Adar             | $\sum_{z \in \Gamma(x) \cap \Gamma(y)} \frac{1}{\log  \Gamma(z) }$  | second |
| resource allocation     | $\sum_{z \in \Gamma(x) \cap \Gamma(y)} \frac{1}{ \Gamma(z) }$   | second |
| Katz                    | $\sum_{l=1}^{\infty} \beta^l  \text{walks}^{\langle l \rangle}(x, y) $  | high   |
| PageRank                | $[\pi_x]_y + [\pi_y]_x$   | high   |
| SimRank                 | $\gamma \frac{\sum_{a \in \Gamma(x)} \sum_{b \in \Gamma(y)} \text{score}(a, b)}{ \Gamma(x)  \cdot  \Gamma(y) }$ | high   |

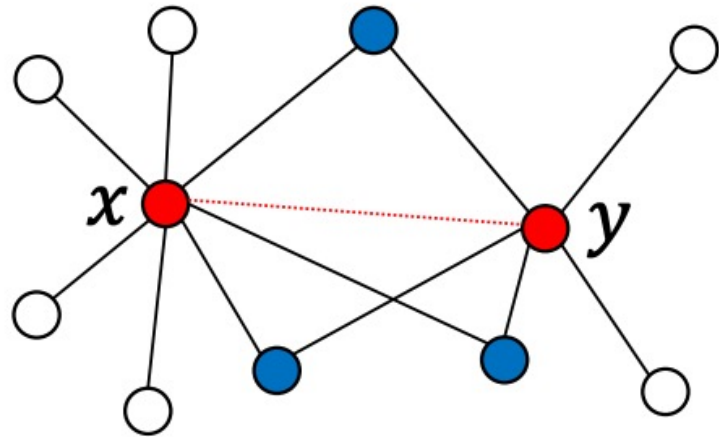
### Pros:

- Easy calculation
- Interpretable
- Not dependent on training

# Subgraph Sampling for Community Information



Heuristic methods for link prediction: **Common Neighbor**



$\Gamma(x)$  denotes the neighbor set of node  $x$

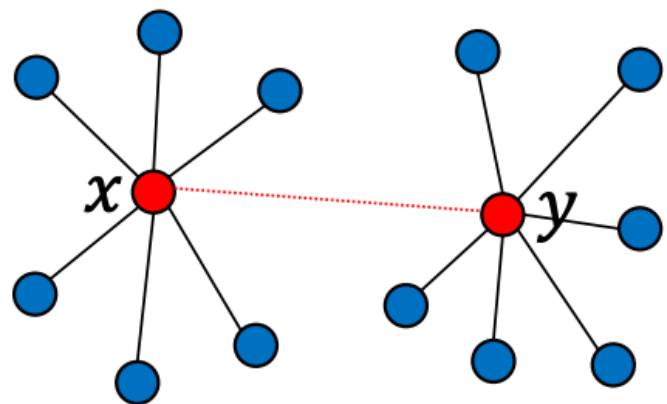
$x$  and  $y$  are likely to have a link if they have many common neighbors

common neighbors (CN):  
 $|\Gamma(x) \cap \Gamma(y)|$

# Subgraph Sampling for Community Information



Heuristic methods for link prediction: **Preferential Attachment**



$\Gamma(x)$  denotes the neighbor set of node  $x$

$x$  is likely to connect to  $y$  if  $y$  has many connections

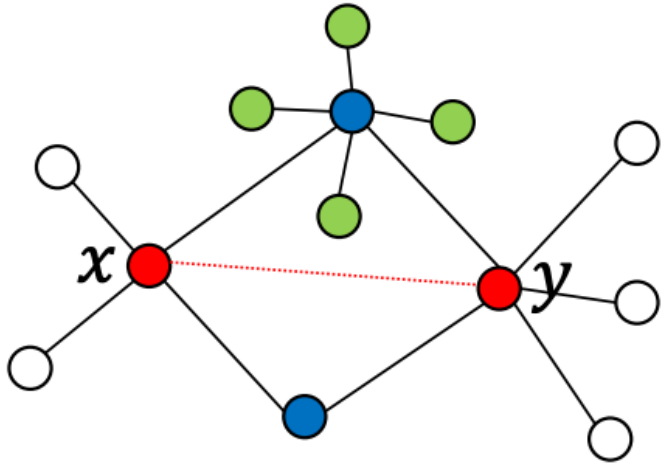
preferential attachment (PA):  
 $|\Gamma(x)| \cdot |\Gamma(y)|$



# Subgraph Sampling for Community Information



Heuristic methods for link prediction: **Adamic-Adar**



Adamic-Adar (AA):

$$\sum_{z \in \Gamma(x) \cap \Gamma(y)} \frac{1}{\log |\Gamma(z)|}$$

$\Gamma(x)$  denotes the neighbor set of node  $x$

- Weighted common neighbors
- Popular common neighbors contribute less

# Subgraph Sampling for Community Information



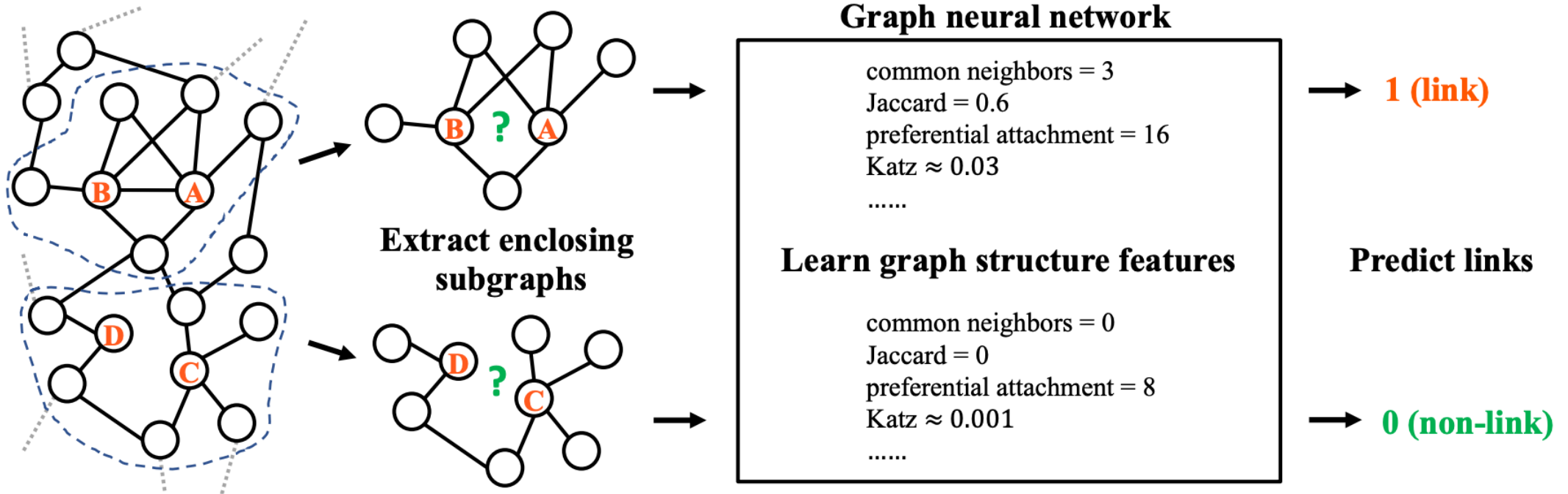
Heuristic methods for link prediction: **Problems**

- Rule-based strategy, not learning-based
- Strong assumption on link formation mechanisms, which might only work well on certain graphs



# Subgraph Sampling for Community Information

Sample a subgraph to enhance link prediction: **SEAL**



# Subgraph Sampling for Community Information

Sample a subgraph to enhance link prediction: **SEAL**



## Comparison with heuristic methods (AUC)

| Data   | CN         | Jaccard    | PA         | AA         | RA         | Katz       | PR                | SR         | ENS        | WLK               | WLNK       | SEAL              |
|--------|------------|------------|------------|------------|------------|------------|-------------------|------------|------------|-------------------|------------|-------------------|
| USAir  | 93.80±1.22 | 89.79±1.61 | 88.84±1.45 | 95.06±1.03 | 95.77±0.92 | 92.88±1.42 | 94.67±1.08        | 78.89±2.31 | 88.96±1.44 | <b>96.63±0.73</b> | 95.95±1.10 | <b>96.62±0.72</b> |
| NS     | 94.42±0.95 | 94.43±0.93 | 68.65±2.03 | 94.45±0.93 | 94.45±0.93 | 94.85±1.10 | 94.89±1.08        | 94.79±1.08 | 97.64±0.25 | 98.57±0.51        | 98.61±0.49 | <b>98.85±0.47</b> |
| PB     | 92.04±0.35 | 87.41±0.39 | 90.14±0.45 | 92.36±0.34 | 92.46±0.37 | 92.92±0.35 | 93.54±0.41        | 77.08±0.80 | 90.15±0.45 | 93.83±0.59        | 93.49±0.47 | <b>94.72±0.46</b> |
| Yeast  | 89.37±0.61 | 89.32±0.60 | 82.20±1.02 | 89.43±0.62 | 89.45±0.62 | 92.24±0.61 | 92.76±0.55        | 91.49±0.57 | 82.36±1.02 | 95.86±0.54        | 95.62±0.52 | <b>97.91±0.52</b> |
| C.ele  | 85.13±1.61 | 80.19±1.64 | 74.79±2.04 | 86.95±1.40 | 87.49±1.41 | 86.34±1.89 | <b>90.32±1.49</b> | 77.07±2.00 | 74.94±2.04 | 89.72±1.67        | 86.18±1.72 | <b>90.30±1.35</b> |
| Power  | 58.80±0.88 | 58.79±0.88 | 44.33±1.02 | 58.79±0.88 | 58.79±0.88 | 65.39±1.59 | 66.00±1.59        | 76.15±1.06 | 79.52±1.78 | 82.41±3.43        | 84.76±0.98 | <b>87.61±1.57</b> |
| Router | 56.43±0.52 | 56.40±0.52 | 47.58±1.47 | 56.43±0.51 | 56.43±0.51 | 38.62±1.35 | 38.76±1.39        | 37.40±1.27 | 47.58±1.48 | 87.42±2.08        | 94.41±0.88 | <b>96.38±1.45</b> |
| E.coli | 93.71±0.39 | 81.31±0.61 | 91.82±0.58 | 95.36±0.34 | 95.95±0.35 | 93.50±0.44 | 95.57±0.44        | 62.49±1.43 | 91.89±0.58 | 96.94±0.29        | 97.21±0.27 | <b>97.64±0.22</b> |

SEAL outperforms heuristic methods

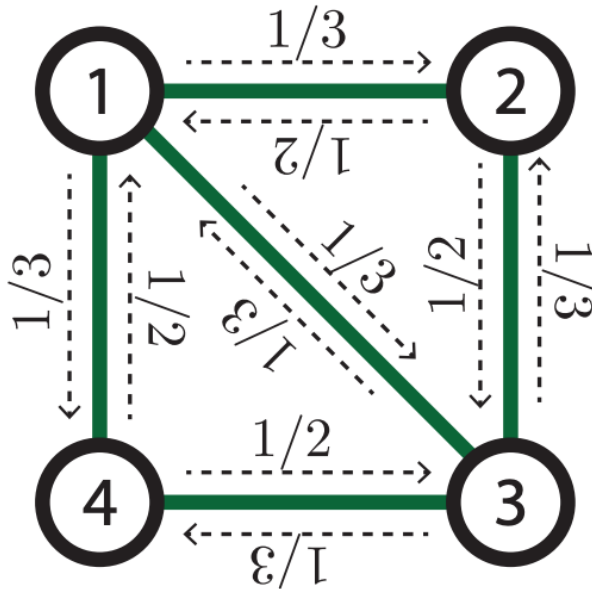
# Subgraph Sampling for Community Information



One step further: **WalkPool**

SEAL average the node embedding in subgraph for link prediction

**WalkPool** proposes that average pooling is suboptimal. Instead, using random walk to extract structural information

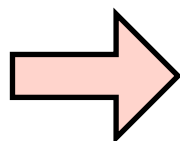
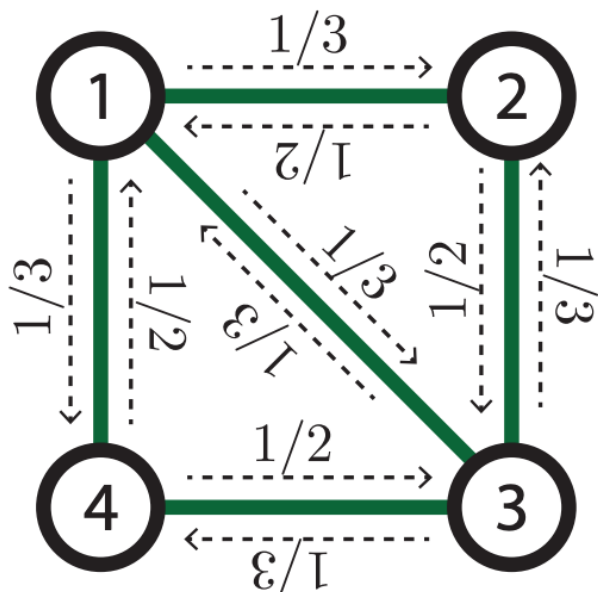


$$\mathbf{P} = \begin{pmatrix} 0 & 1/3 & 1/3 & 1/3 \\ 1/2 & 0 & 1/2 & 0 \\ 1/3 & 1/3 & 0 & 1/3 \\ 1/2 & 0 & 1/2 & 0 \end{pmatrix}$$

Random walk probabilities

# Subgraph Sampling for Community Information

One step further: **WalkPool**



| Features   | $\tau = 2$ walk length   |
|--|--|
| node <sub>1</sub> <sup><math>\tau</math></sup>   | $\left( \frac{1}{3} \cdot \frac{1}{2} + \frac{1}{3} \cdot \frac{1}{2} + \frac{1}{3} \cdot \frac{1}{3} \right)$ <div style="display: flex; justify-content: space-around; align-items: center;"> <div style="text-align: center;"> <math>\textcircled{1} \rightarrow \textcircled{2} \rightarrow \textcircled{1}</math> </div> <div style="text-align: center;"> <math>\textcircled{1} \rightarrow \textcircled{4} \rightarrow \textcircled{1}</math> </div> <div style="text-align: center;"> <math>\textcircled{1} \rightarrow \textcircled{3} \rightarrow \textcircled{1}</math> </div> </div> |
| link <sub>1,2</sub> <sup><math>\tau</math></sup> | $\left( \frac{1}{3} \cdot \frac{1}{2} \right)$ <div style="text-align: center;"> <math>\textcircled{1} \rightarrow \textcircled{3} \rightarrow \textcircled{2}</math> </div>   |
| graph <sup><math>\tau</math></sup>               | $\text{node}_1^2 + \text{node}_2^2 + \text{node}_3^2 + \text{node}_4^2$  |

Use these random walk features for link prediction via an MLP

# Graph Learning Enhanced by Knowledge from Data

## Chapter summary



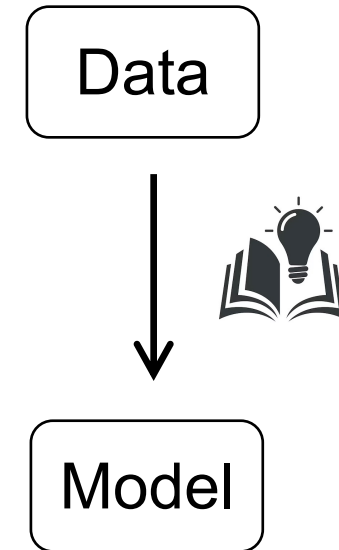
Knowledge can be obtained from

### 1) **single-instance** level perception

- Node sampling for node positions


### 2) **multiple-instance** level perception

- Path sampling for positional and semantic information
- Subgraph sampling for community information



# Tutorial Outline



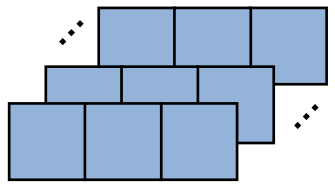
- Preliminaries and Foundations
- Graph Learning Enhanced by Knowledge from Data
-  • **Graph Learning Enhanced by Knowledge from Models**
- Graph Learning Enhanced by Knowledge from Humans and Domains
- Graph Learning Enhanced by Knowledge from External Sources
- Knowledge-enhanced Graph Learning for Real-world Applications
- Summary and Future Directions





# Models Contain Rich Knowledge

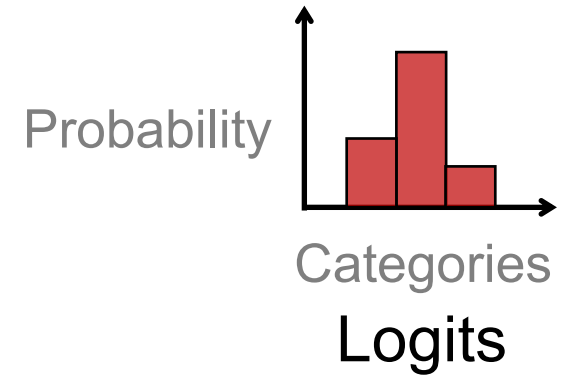
## What is knowledge in models?



Learned Embeddings

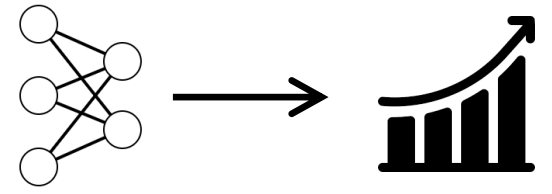
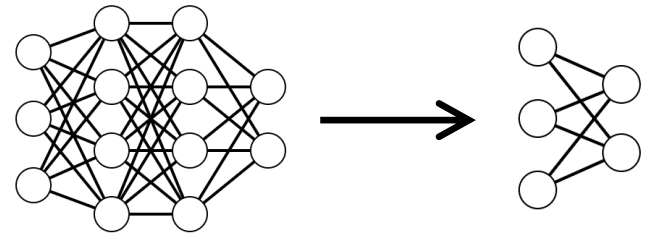
|    | p1 | p2 | p3 |
|----|----|----|----|
| p1 | 1  | .3 | .7 |
| p2 | .3 | 1  | .2 |
| p3 | .7 | .2 | 1  |

Data Similarity



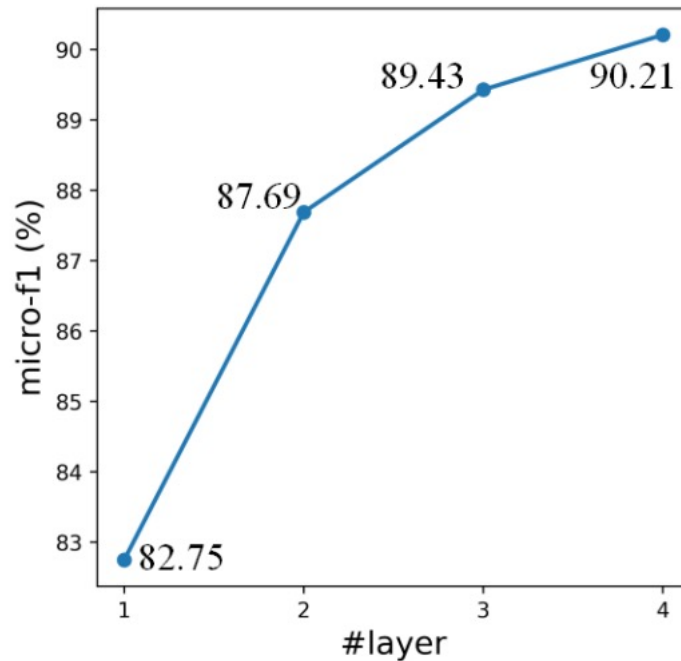
## Why we need to emphasize model knowledge?

- Two goals:
- Model Compression
  - Performance Improvement

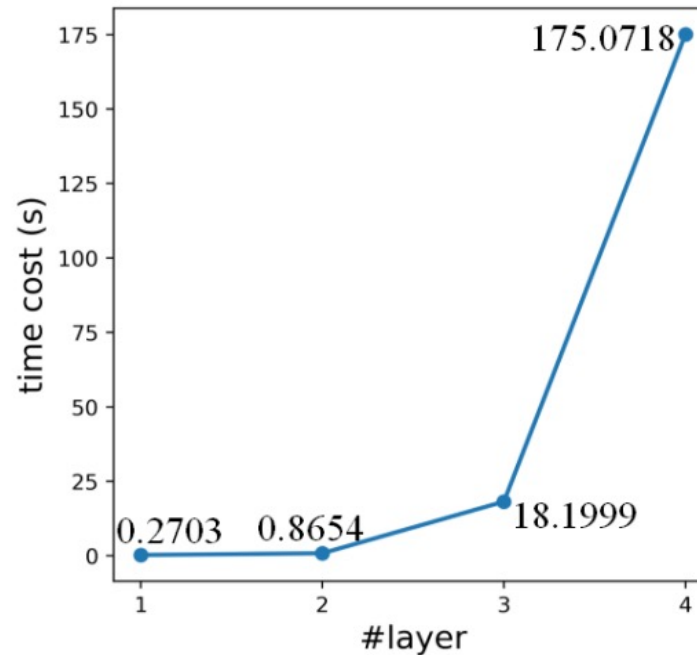


# Models Contain Rich Knowledge

Deep GNNs are expensive



Micro-f1 of GAT on Facebook



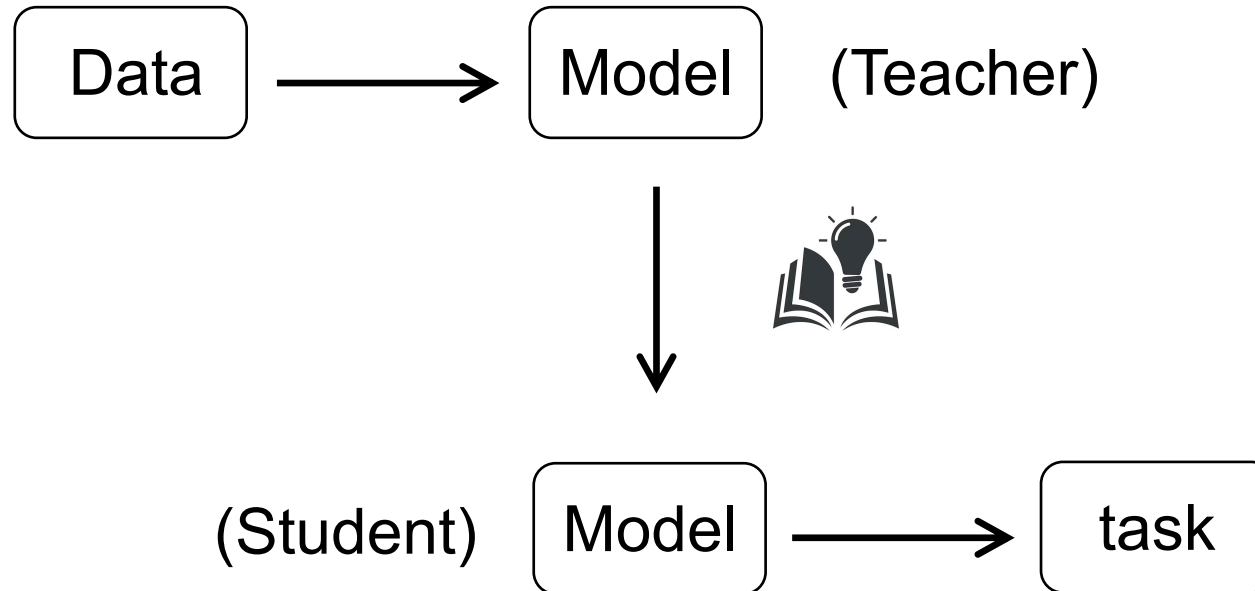
Time cost of GAT on Facebook

The more layers,  
the higher performance,  
the more expensive

# Models Contain Rich Knowledge

How can we emphasize model knowledge?

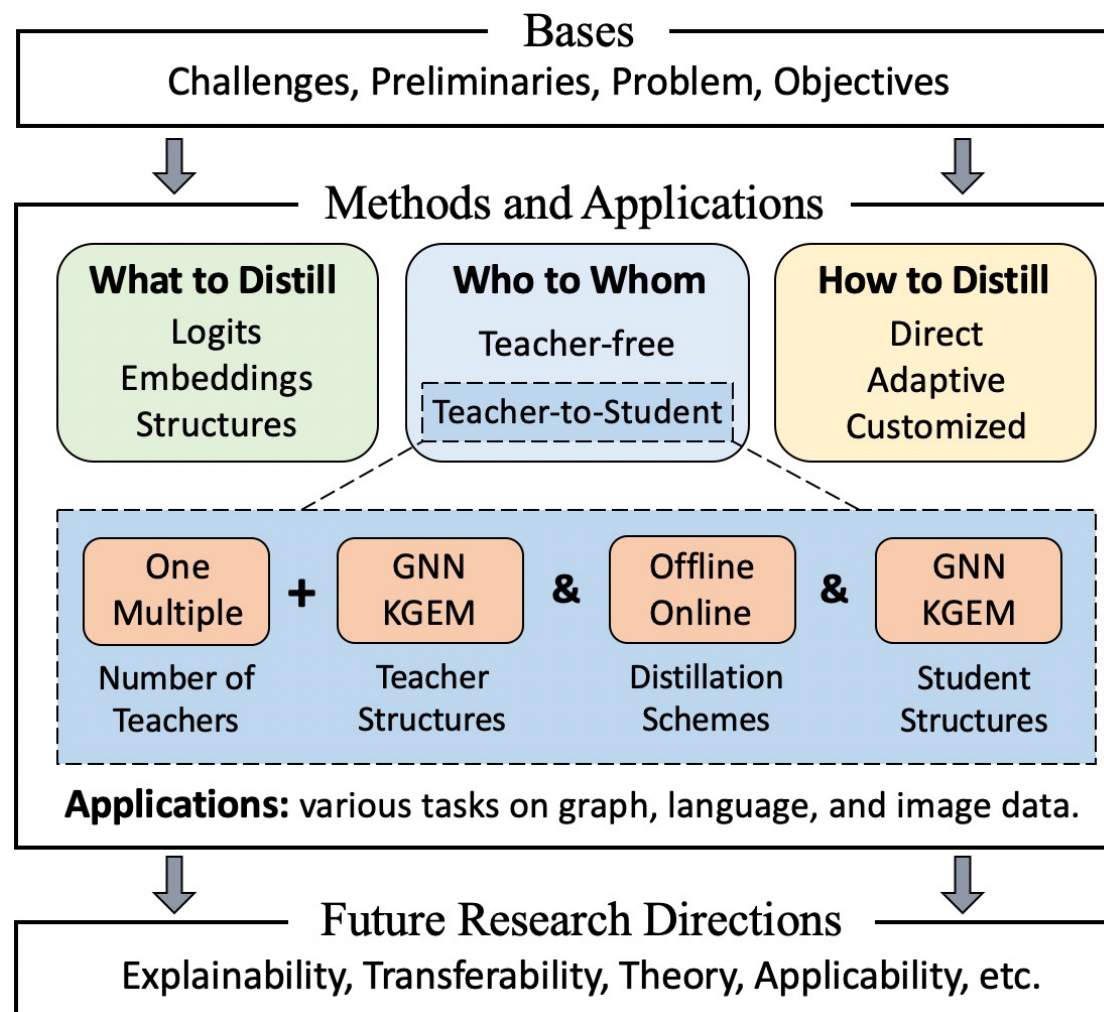
Distilling knowledge from models!



# A Systematic Framework

## Three Main Aspects:

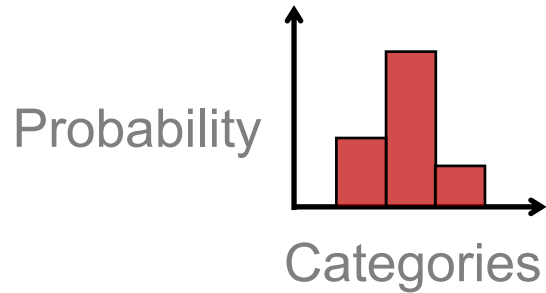
- What to Distill
- Who to Whom
- How to Distill



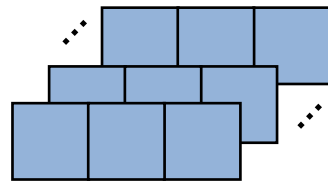
# A Systematic Framework



## What to distill?



**Logits** denote the inputs to the final SoftMax function and represent the soft label prediction



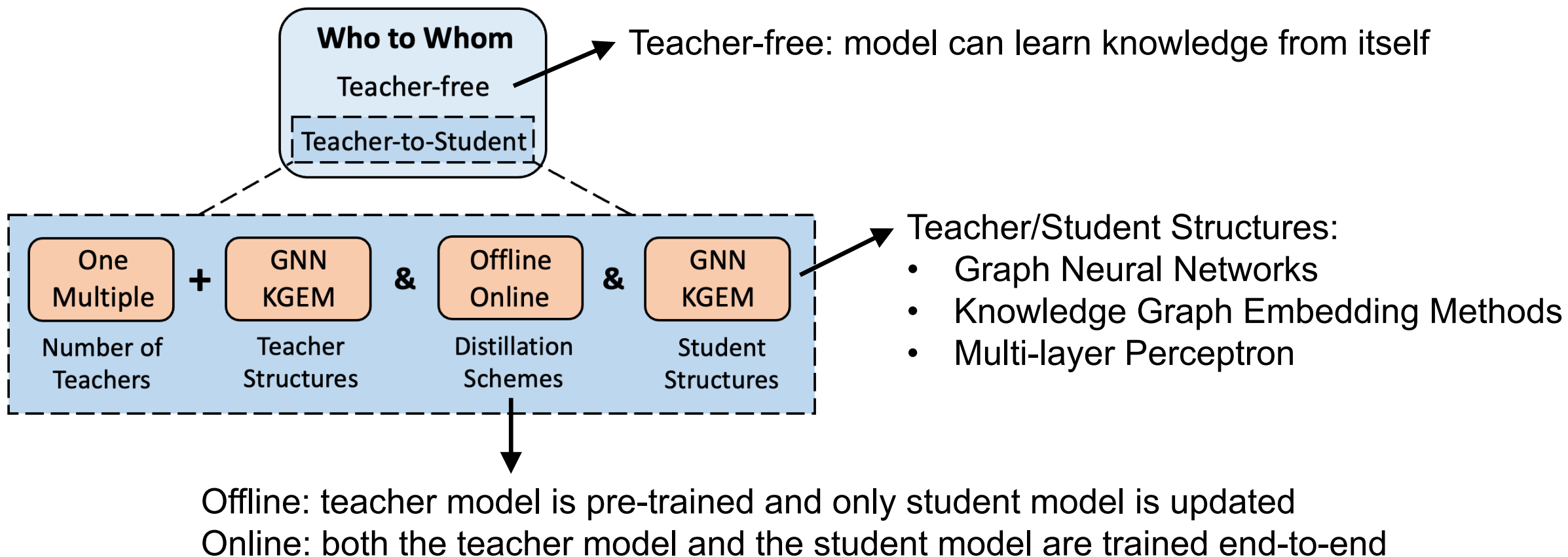
**Embeddings** are learned node embeddings from the intermediate layers of teacher models

|    |    |    |
|----|----|----|
| 1  | .3 | .7 |
| .3 | 1  | .2 |
| .7 | .2 | 1  |

**Structures** depict the connectivity and relationships between the elements in graph

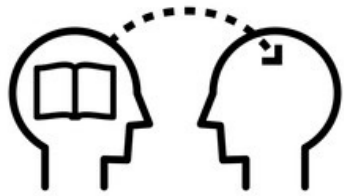
# A Systematic Framework

## Who to whom?

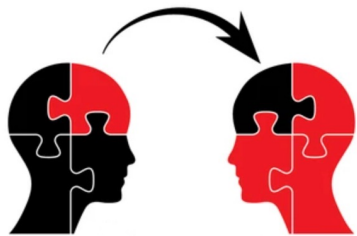


# A Systematic Framework

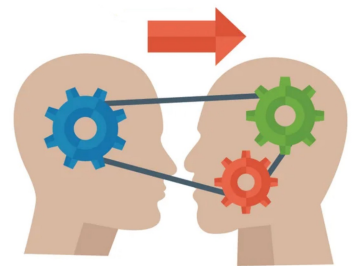
## How to distill?



**Direct distillation** minimizes the divergences between the knowledge of the teacher and the student directly



**Adaptive distillation** considers the significance of teacher knowledge adaptively



**Customized distillation** enables diversified designs according to the distinct objective of tasks in different scenarios





# Covered Topics for Knowledge from Models



## How to obtain knowledge from

- 1) Logits
- 2) Embeddings
- 3) Structures

## A recent learning strategy

- 4) Distilling an MLP to replace GNNs

# Covered Topics for Knowledge from Models



How to obtain knowledge from



1) Logits

2) Embeddings

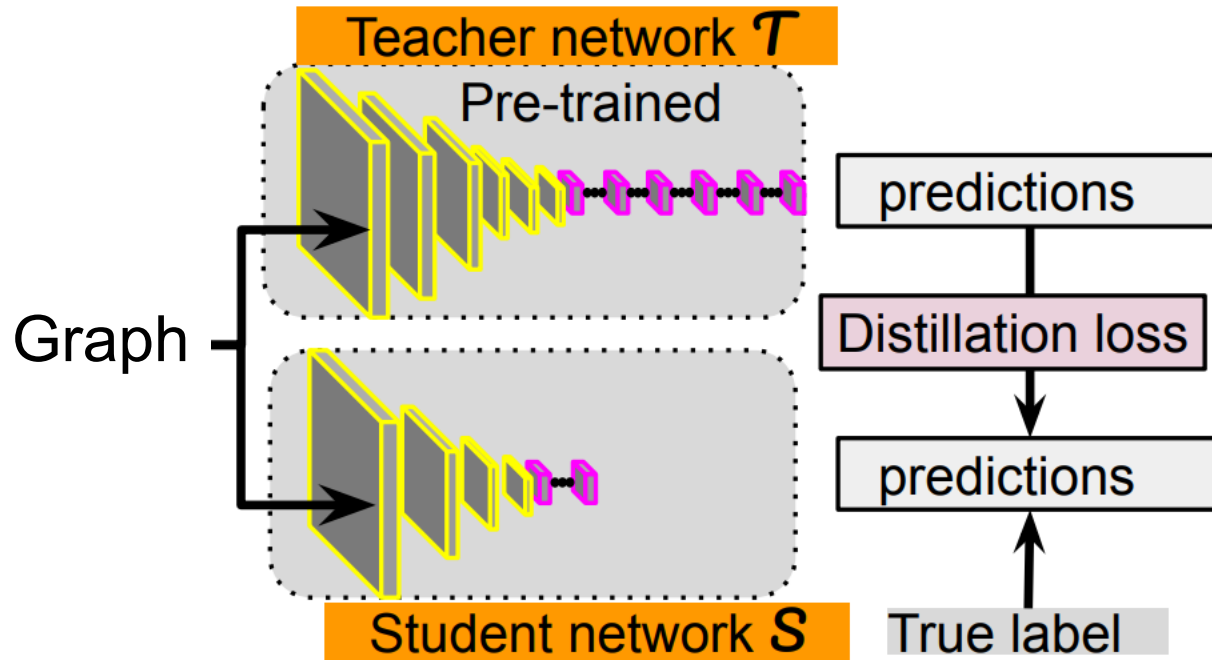
3) Structures

A recent learning strategy

4) Distilling an MLP to replace GNNs

# Knowledge Distillation from Logits

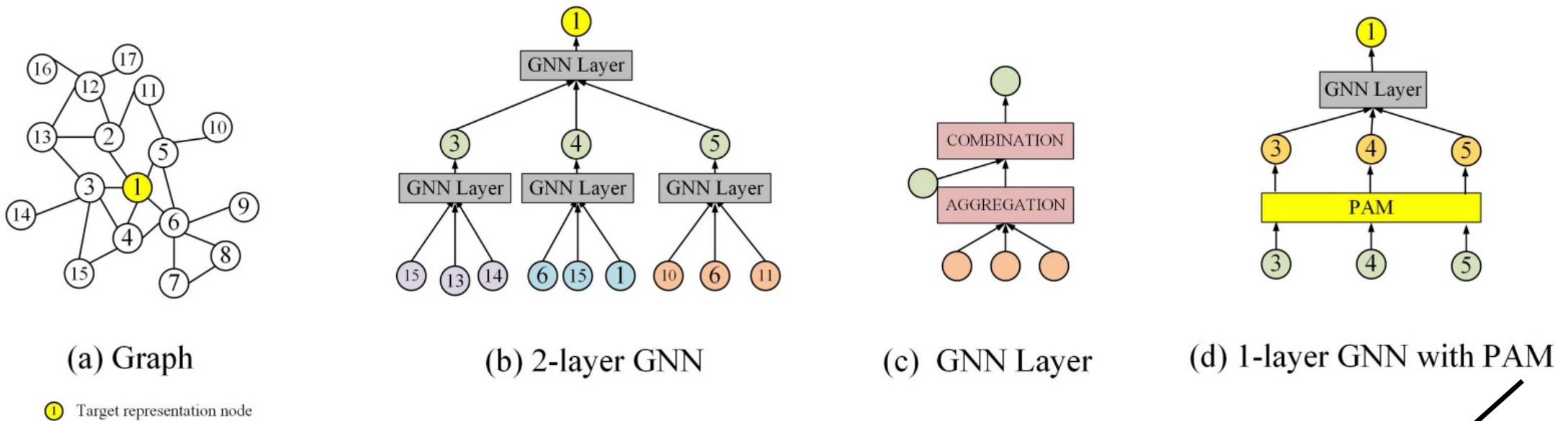
## Illustration



Learning a smaller student network by mimicking the predictions of the teacher

# Knowledge Distillation from Logits

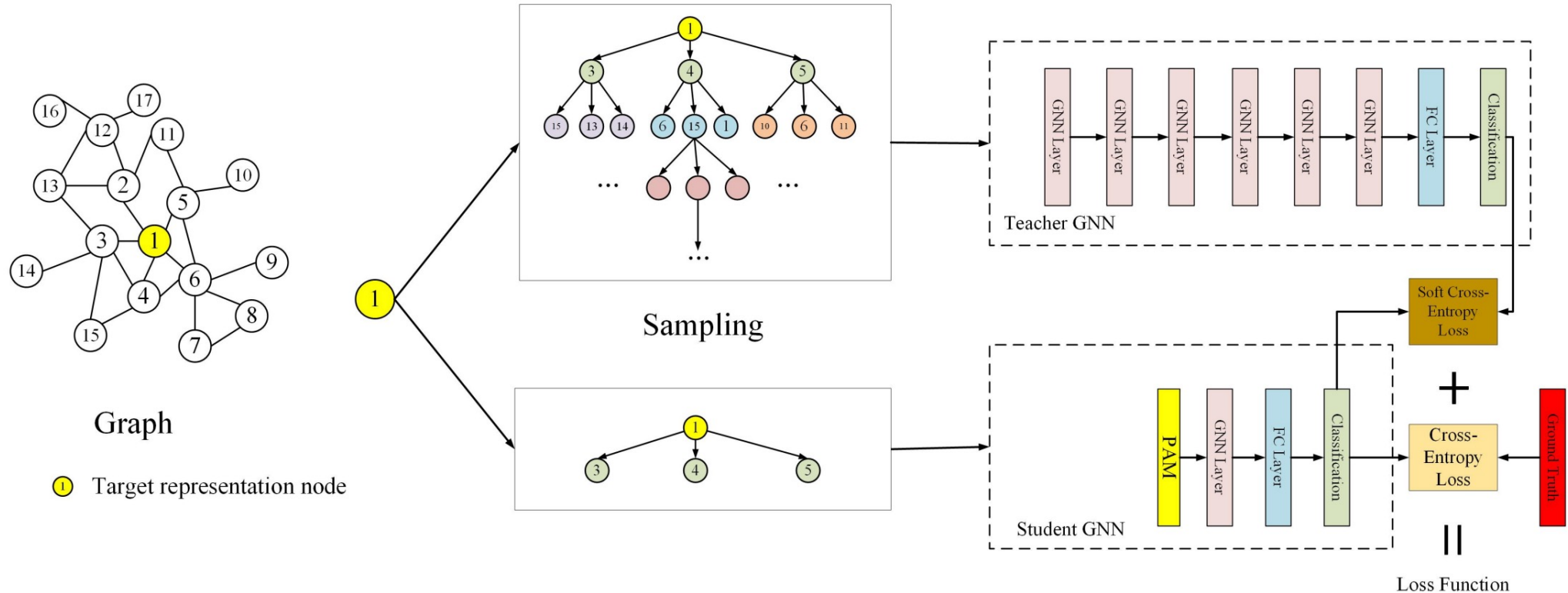
Learn a small GNN with local structure enhanced: TinyGNN



Peer-Aware Module captures the local structure

# Knowledge Distillation from Logits

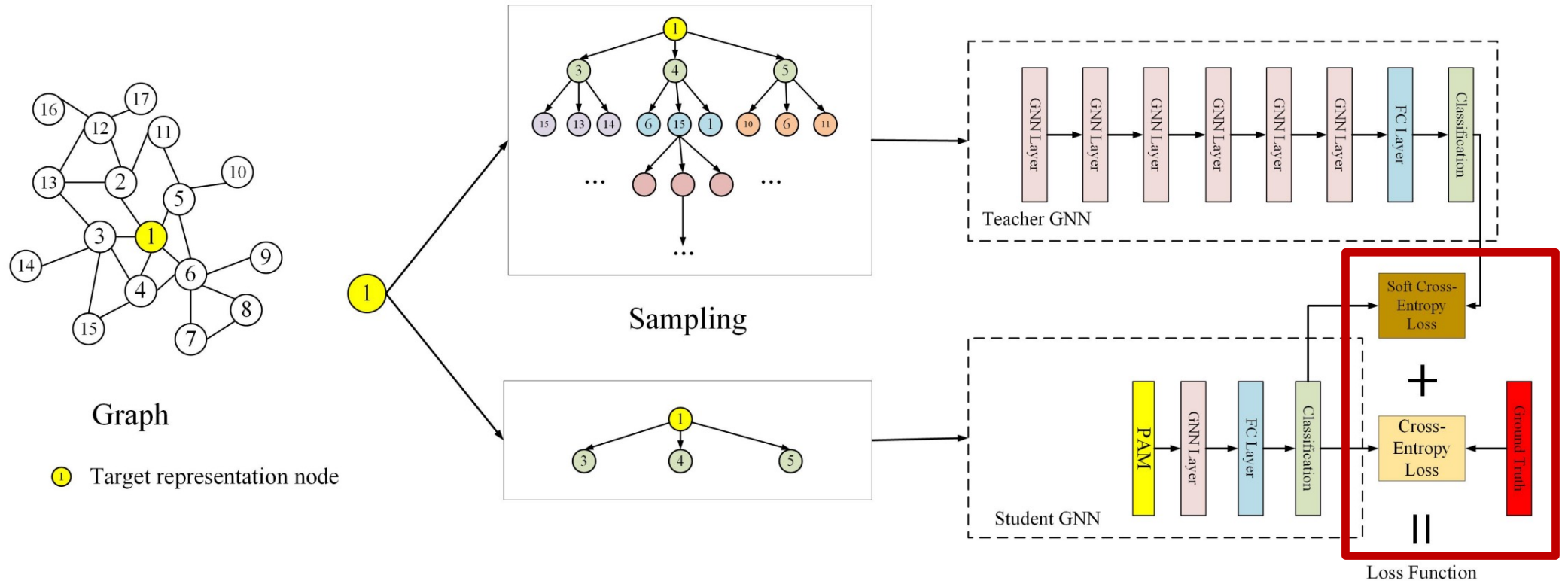
## Learn a small GNN with local structure enhanced: TinyGNN



- Nodes 3, 4, 5 are peer nodes to each other
- Calculate attention between peer nodes to encode local structure

# Knowledge Distillation from Logits

## Learn a small GNN with local structure enhanced: TinyGNN



Learning knowledge from teacher logits and ground truth labels

# Knowledge Distillation from Logits

## Learn a small GNN with local structure enhanced: TinyGNN

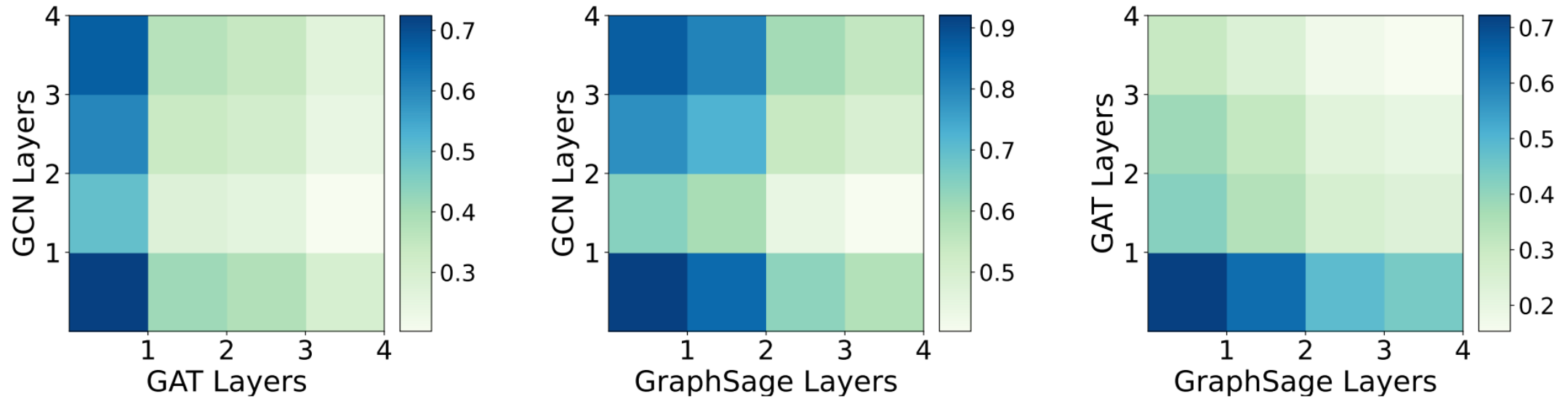


| Model                             | Facebook     |              | Chameleon    |              | Squirrel     |              | AliGraph     |              |
|-----------------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
|                                   | Micro-f1(%)  | Macro-f1(%)  | Micro-f1(%)  | Macro-f1(%)  | Micro-f1(%)  | Macro-f1(%)  | Micro-f1(%)  | Macro-f1(%)  |
| <i>GNN</i> <sub>3</sub> (Teacher) | 89.43        | 88.69        | 39.12        | 38.16        | 30.68        | 30.41        | 60.34        | 53.75        |
| <i>GNN</i> <sub>2</sub> (Base)    | 87.69        | 86.83        | 38.73        | 36.92        | 30.47        | 30.08        | 49.98        | 43.45        |
| <i>GNN</i> <sub>1</sub> (Base)    | 82.75        | 81.65        | 36.88        | 35.79        | 29.83        | 29.24        | 33.77        | 25.22        |
| <i>GNN</i> <sub>1</sub> -NDS      | 83.32        | 82.11        | 37.46        | 36.42        | 30.73        | 30.38        | 35.30        | 24.88        |
| <i>GNN</i> <sub>1</sub> -PAM      | 83.15        | 81.79        | 38.05        | 36.56        | 31.32        | 31.08        | 40.19        | 32.05        |
| <i>TinyGNN</i> <sub>1</sub>       | <b>84.56</b> | <b>83.41</b> | <b>39.71</b> | <b>38.46</b> | <b>32.09</b> | <b>32.00</b> | <b>42.34</b> | <b>32.37</b> |
| <i>GNN</i> <sub>2</sub> -NDS      | 88.90        | 88.15        | 39.51        | 38.21        | 30.98        | 30.52        | 54.72        | 46.97        |
| <i>GNN</i> <sub>2</sub> -PAM      | 88.56        | 87.72        | 40.98        | 39.62        | 32.69        | 31.95        | 57.81        | 50.40        |
| <i>TinyGNN</i> <sub>2</sub>       | <b>89.40</b> | <b>88.66</b> | <b>41.17</b> | <b>40.01</b> | <b>33.33</b> | <b>33.17</b> | <b>61.12</b> | <b>54.17</b> |

- 1 layer TinyGNN performs poorly
- 2 layers TinyGNN can outperform the 3 layers GNN teacher

# Knowledge Distillation from Logits

Different GNNs learn different knowledge

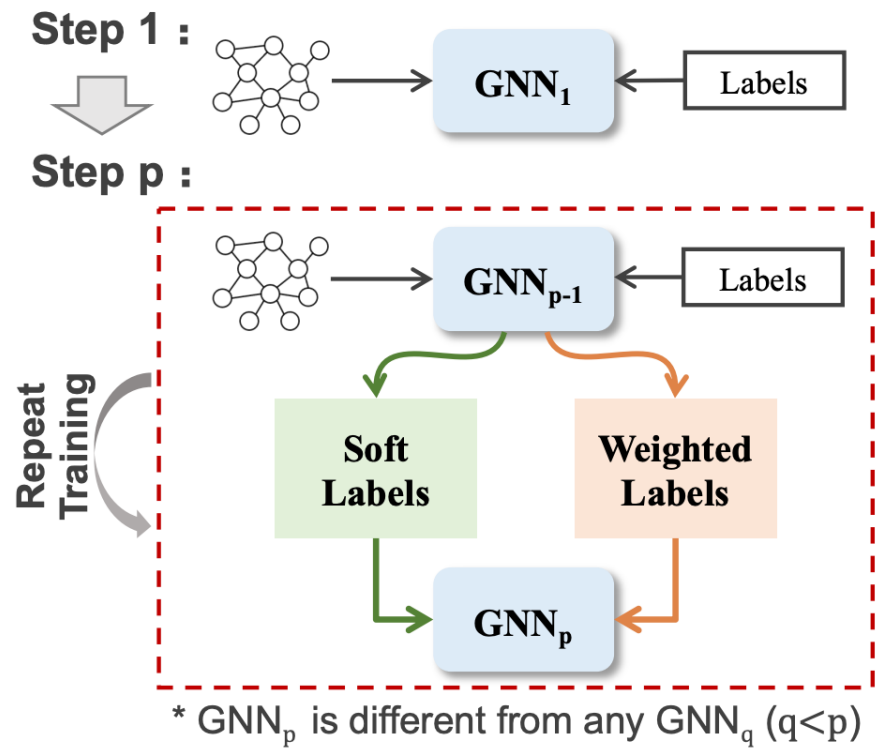


- GCN aggregates neighborhoods with predefined weights
- GAT aggregates neighborhoods using learnable weights
- GraphSage randomly samples neighbors during aggregation

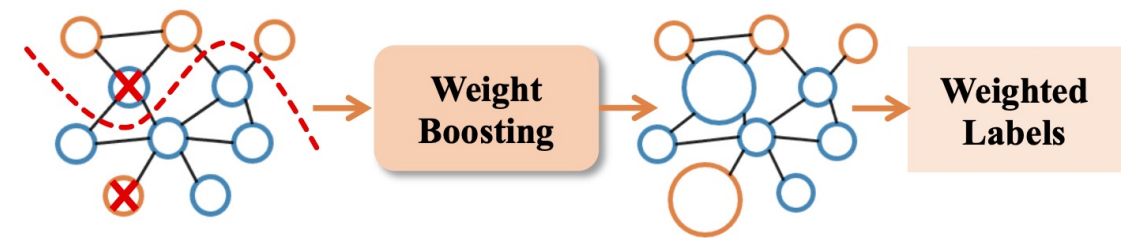


# Knowledge Distillation from Logits

## Learn a GNN from multiple GNNs: BGNN



- Distill several times, once a different GNN
- Weighted Labels encourage student to learn adaptively, by focusing more on the misclassified samples



# Knowledge Distillation from Logits

## Learn a GNN from multiple GNNs: BGNN



| Student | Method                       | Graph Classification |                   |                   | Node Classification |                   |                   |                   |
|---------|------------------------------|----------------------|-------------------|-------------------|---------------------|-------------------|-------------------|-------------------|
|         |                              | COLLAB               | IMDB              | ENZYMES           | CORA                | CITeseer          | PUBMED            | A-COMPUTERS       |
| GCN     | BAN (Furlanello et al. 2018) | 81.60±0.40           | 78.50±1.00        | 66.67±1.67        | 83.17±0.26          | 72.47±0.21        | 79.87±0.21        | 88.78±0.05        |
|         | MulDE (Wang et al. 2021)     | 80.86±1.18           | 77.50±1.20        | 67.33±0.90        | 82.53±0.05          | 72.33±0.09        | 78.73±0.09        | 88.41±0.12        |
|         | BGNN(s)                      | 82.73±0.34           | 79.33±0.47        | 69.44±0.79        | 83.97±0.17          | 73.87±0.24        | 80.73±0.25        | <b>89.58±0.03</b> |
|         | BGNN(m)-ST                   | 82.87±0.09           | <b>79.67±0.27</b> | <b>71.12±2.45</b> | <b>84.83±0.25</b>   | 73.60±0.14        | 80.20±0.08        | 89.03±0.02        |
|         | BGNN(m)-TS                   | <b>83.40±0.15</b>    | 79.00±0.00        | 70.00±1.36        | 84.40±0.22          | <b>74.87±0.25</b> | <b>80.90±0.00</b> | 89.02±0.03        |

BGNN outperforms existing methods on multi-teacher KD

# Covered Topics for Knowledge from Models



How to obtain knowledge from

1) Logits



2) Embeddings

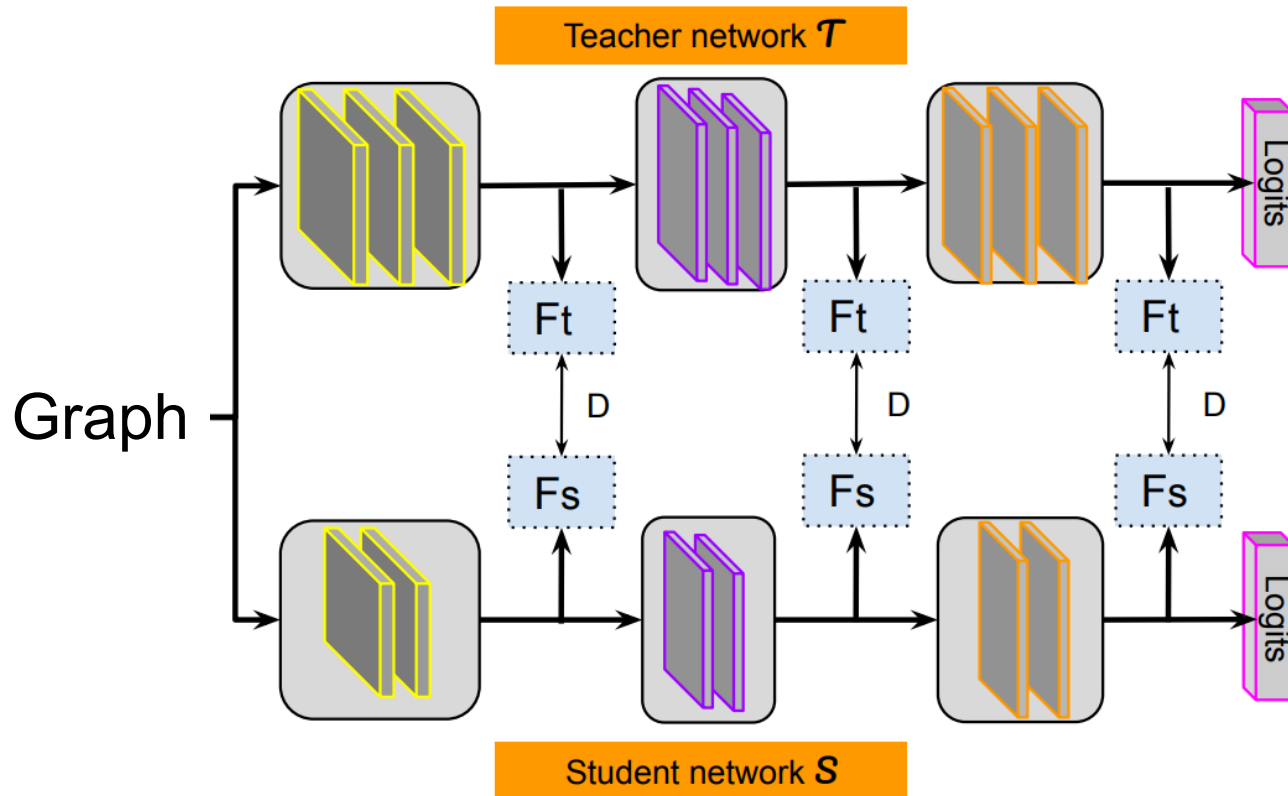
3) Structures

A recent learning strategy

4) distilling an MLP to replace GNNs

# Knowledge Distillation from Embeddings

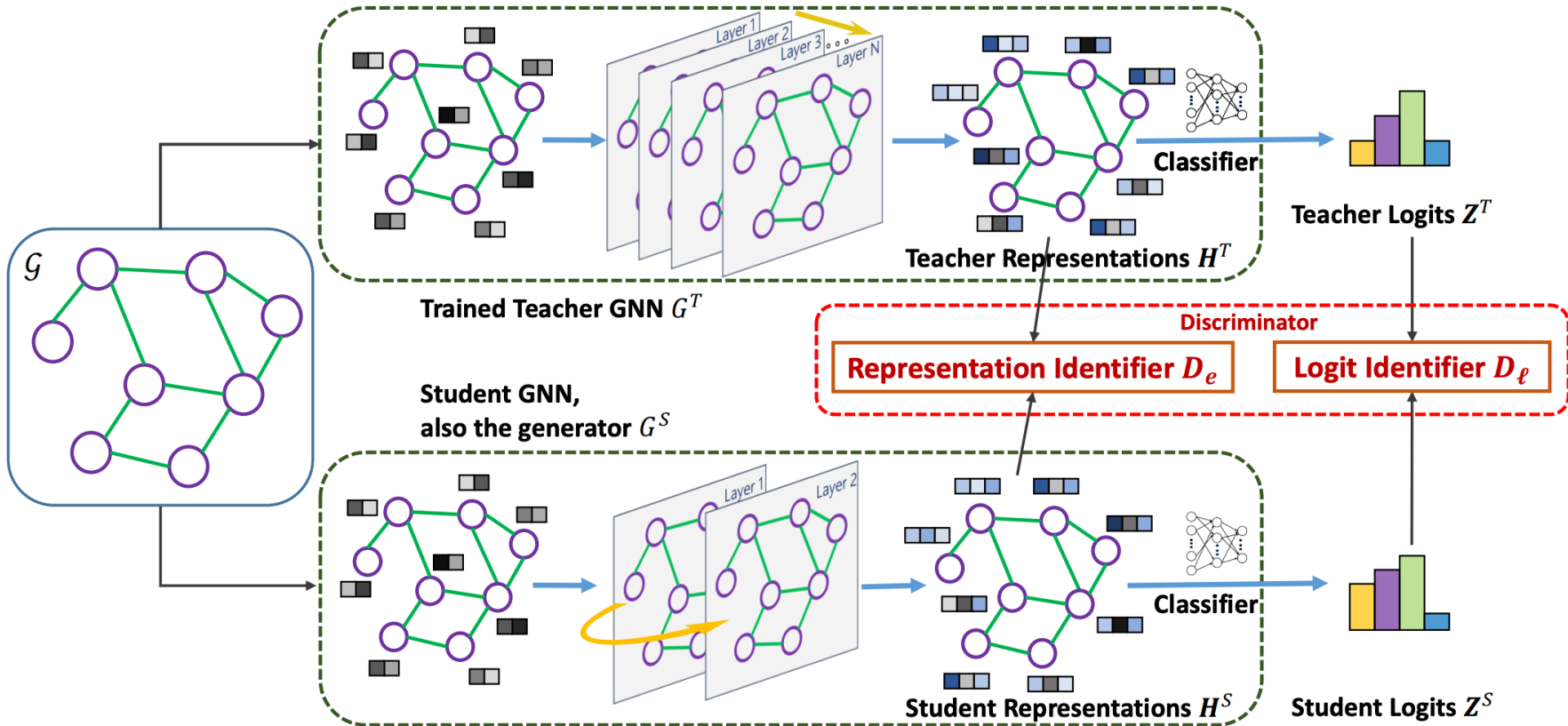
## Illustration



- Beyond learning teacher's logits, the student should also **learn teacher's embeddings**
- Logits contain inter-class correlations
- Embeddings contain **inter-node correlations**

# Knowledge Distillation from Embeddings

Enable the student to learn embeddings: GraphAKD



- **GNN Generator:** produce similar embeddings and logits
- **Discriminator:** distinguish the output of teacher and student

# Knowledge Distillation from Embeddings

Enable the student to learn embeddings: GraphAKD



| Datasets | Teacher |       | Vanilla Student |             |              | Student trained with GraphAKD |             |           |                 |               |
|----------|---------|-------|-----------------|-------------|--------------|-------------------------------|-------------|-----------|-----------------|---------------|
|          | Model   | Perf. | #Params         | Model       | O. Perf.     | R. Perf.                      | Perf.       | #Params   | Perf. Impv. (%) | #Params Decr. |
| Cora     | GCNII   | 85.5  | 616,519         | GCN         | <u>81.5</u>  | 78.3 ±0.9                     | 83.6 ±0.8   | 96,633    | 2.1             | 84.3%         |
| CiteSeer | GCNII   | 73.4  | 5,144,070       | GCN         | <u>71.1</u>  | 68.6 ±1.1                     | 72.9 ±0.4   | 1,016,156 | 1.8             | 80.2%         |
| PubMed   | GCNII   | 80.3  | 1,177,603       | GCN         | <u>79.0</u>  | 78.1 ±1.0                     | 81.3 ±0.4   | 195,357   | 2.3             | 83.4%         |
| Flickr   | GCNII   | 56.20 | 1,182,727       | GCN         | 49.20        | <u>49.63</u> ±1.19            | 52.95 ±0.24 | 196,473   | 3.32            | 83.4%         |
| Arxiv    | GCNII   | 72.74 | 2,148,648       | GCN         | <u>71.74</u> | 71.43 ±0.13                   | 73.05 ±0.22 | 242,426   | 1.31            | 88.7%         |
| Reddit   | GCNII   | 96.77 | 691,241         | GCN         | 93.30        | <u>94.12</u> ±0.04            | 95.15 ±0.02 | 234,655   | 1.03            | 66.1%         |
| Yelp     | GCNII   | 65.14 | 2,306,660       | Cluster-GCN | 59.15        | <u>59.63</u> ±0.51            | 60.63 ±0.42 | 431,950   | 1.00            | 81.3%         |
| Products | GAMLP   | 84.59 | 3,335,831       | Cluster-GCN | <u>76.21</u> | 74.99 ±0.76                   | 81.45 ±0.47 | 682,449   | 5.24            | 79.5%         |

- Better performance than vanilla student, but perform poorly than the teacher
- Smaller model size (parameters decreased by ~80%)

# Knowledge Distillation from Embeddings

Enable the student to learn embeddings: GraphAKD



| Datasets | #Params |         | GPU Memory |         | Inference time |         |
|----------|---------|---------|------------|---------|----------------|---------|
|          | Teacher | Student | Teacher    | Student | Teacher        | Student |
| Cora     | 0.6M    | 0.1M    | 0.22G      | 0.03G   | 40.3ms         | 4.1ms   |
| PubMed   | 1.2M    | 0.2M    | 1.23G      | 0.33G   | 57.3ms         | 5.7ms   |
| Flickr   | 1.2M    | 0.2M    | 2.79G      | 1.49G   | 309.7ms        | 11.9ms  |
| Yelp     | 2.3M    | 0.4M    | 6.28G      | 4.73G   | 3.0s           | 1.5s    |
| Products | 3.3M    | 0.7M    | 6.25G      | 6.20G   | 16.1s          | 7.0s    |

Reduced GPU memory; Faster inference time



# Knowledge Distillation from Embeddings

Enable the student to learn embeddings: GraphAKD



Ablation studies on the learning impacts

| Datasets                 | Cora        | PubMed      | Flickr       | Yelp         | Products     | Molhiv       |
|--------------------------|-------------|-------------|--------------|--------------|--------------|--------------|
| Teacher                  | 85.5        | 80.3        | 56.20        | 65.14        | 84.59        | 84.03        |
| Student                  | 81.5        | 79.0        | 49.20        | 59.15        | 76.21        | 75.58        |
| Embeddings<br>Only $D_e$ | 82.9        | 80.6        | 52.20        | 59.63        | 81.13        | 78.28        |
| Logits<br>Only $D_\ell$  | 82.3        | 81.0        | 52.52        | 60.03        | 79.76        | 78.09        |
| GraphAKD                 | <b>83.6</b> | <b>81.3</b> | <b>52.95</b> | <b>60.63</b> | <b>81.45</b> | <b>79.16</b> |

- Learning from embeddings is beneficial
- Combining the learning from embeddings and logits is even better




# Covered Topics for Knowledge from Models



How to obtain knowledge from

1) Logits

2) Embeddings

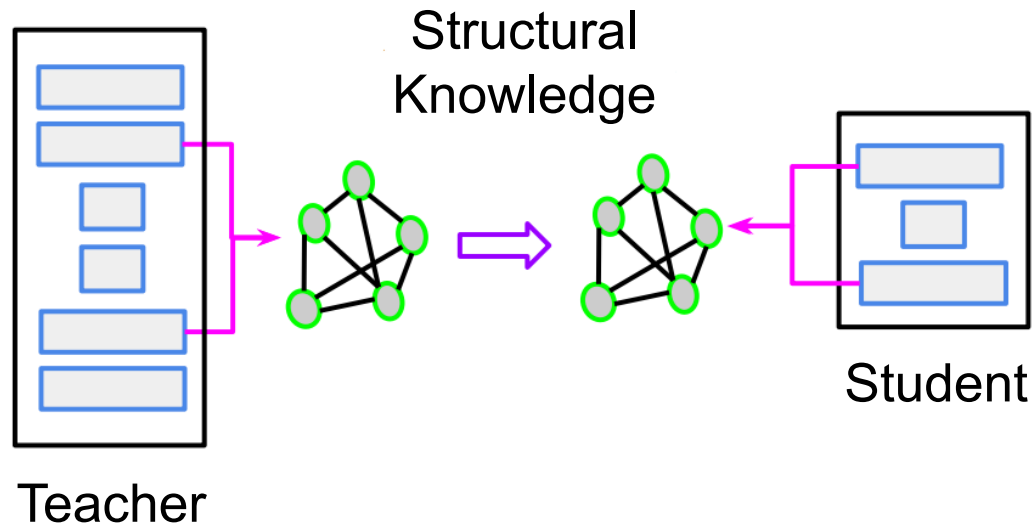
 3) Structures

A recent learning strategy

4) distilling an MLP to replace GNNs

# Knowledge Distillation from Structures

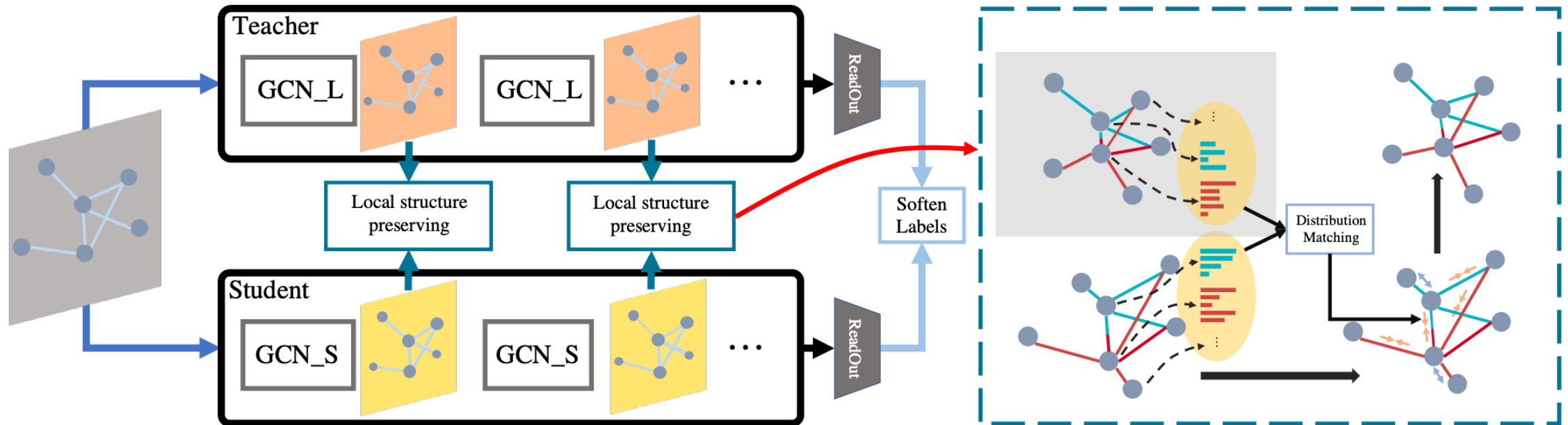
## Illustration



The student should also learn graph structural and topology information from the teacher

# Knowledge Distillation from Structures

Enable the student to capture the local structure: LSP



- 1st step: compute the distribution of the local structure for each node (denoted by node similarity to each other)
- 2nd step: match the distributions of the teacher with that of the student

# Knowledge Distillation from Structures

Enable the student to capture the local structure: LSP



Model details: GAT

| Model   | Layers | Attention heads | Hidden features |
|---------|--------|-----------------|-----------------|
| Teacher | 3      | 4,4,6           | 256,256,121     |
| Student | 5      | 2,2,2,2,2       | 68,68,68,68,121 |

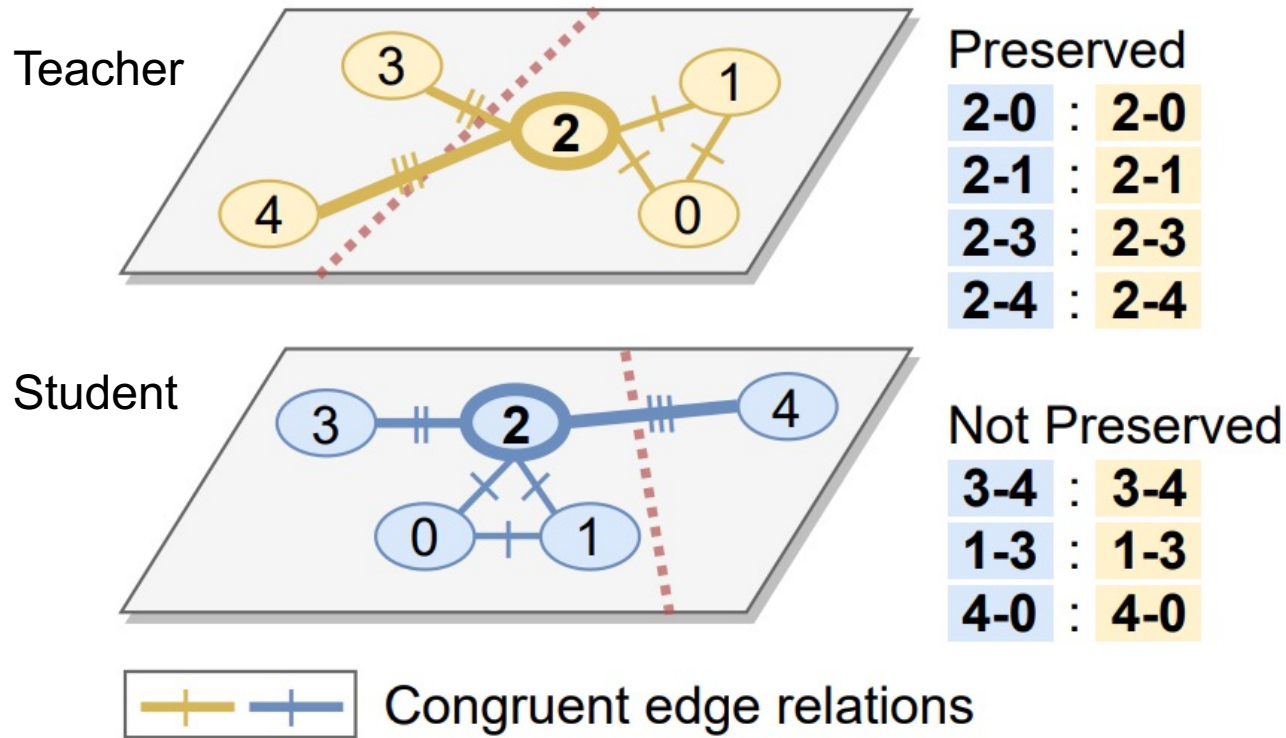
Node classification on PPI

| Model               | Params | RunTime | Training  | F1 Score |
|---------------------|--------|---------|-----------|----------|
| Teacher             | 3.64M  | 48.5ms  | 1.7s/3.4G | 97.6     |
| Student_Full        | 0.16M  | 41.3ms  | 1.3s/1.2G | 95.7     |
| Student_KD [14]     | -      | -       | -         | -        |
| Student_AT [47]     | 0.16M  | 41.3ms  | 1.9s/1.4G | 95.4     |
| Student_FitNet [31] | 0.16M  | 41.3ms  | 2.4s/1.6G | 95.6     |
| Student_LSP (Ours)  | 0.16M  | 41.3ms  | 2.0s/1.5G | 96.1     |

Reduced parameters; Comparable F1 to the teacher

# Knowledge Distillation from Structures

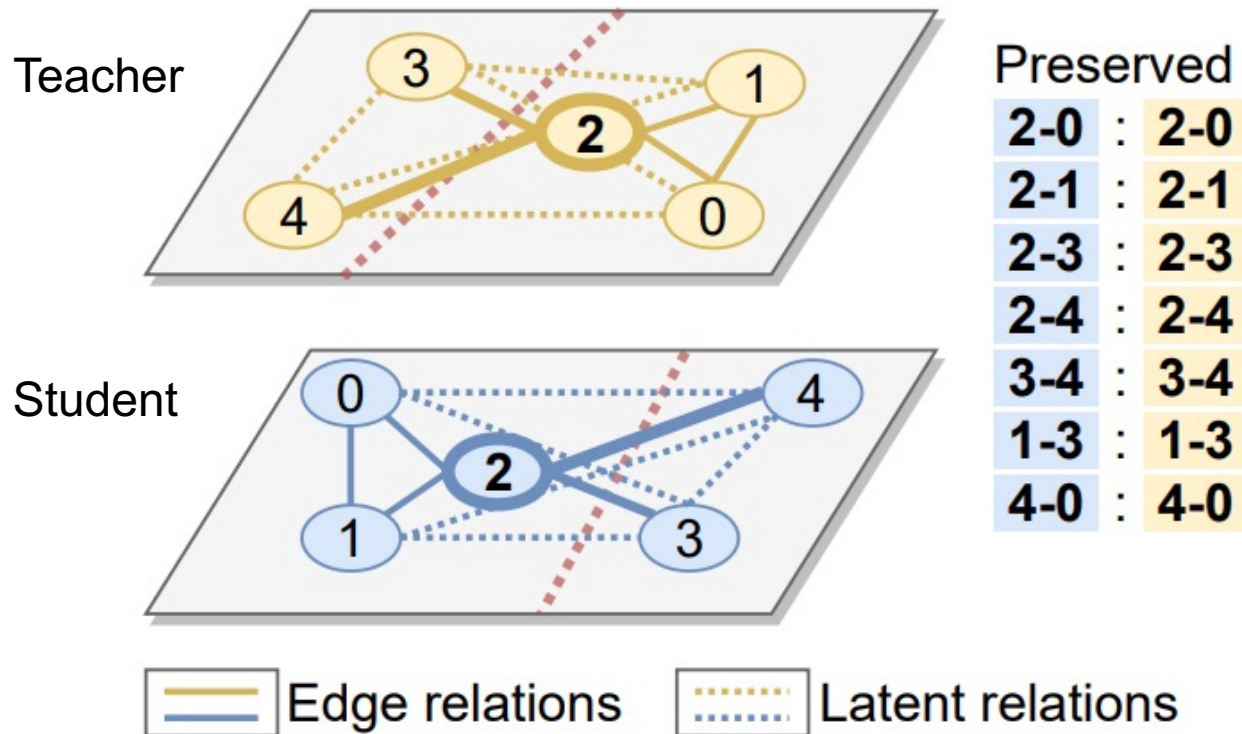
## The strategy of LSP



Considering **pairwise relationships**, but not considering latent interactions among disconnected nodes.

# Knowledge Distillation from Structures

One step further: enhance the student with global structure

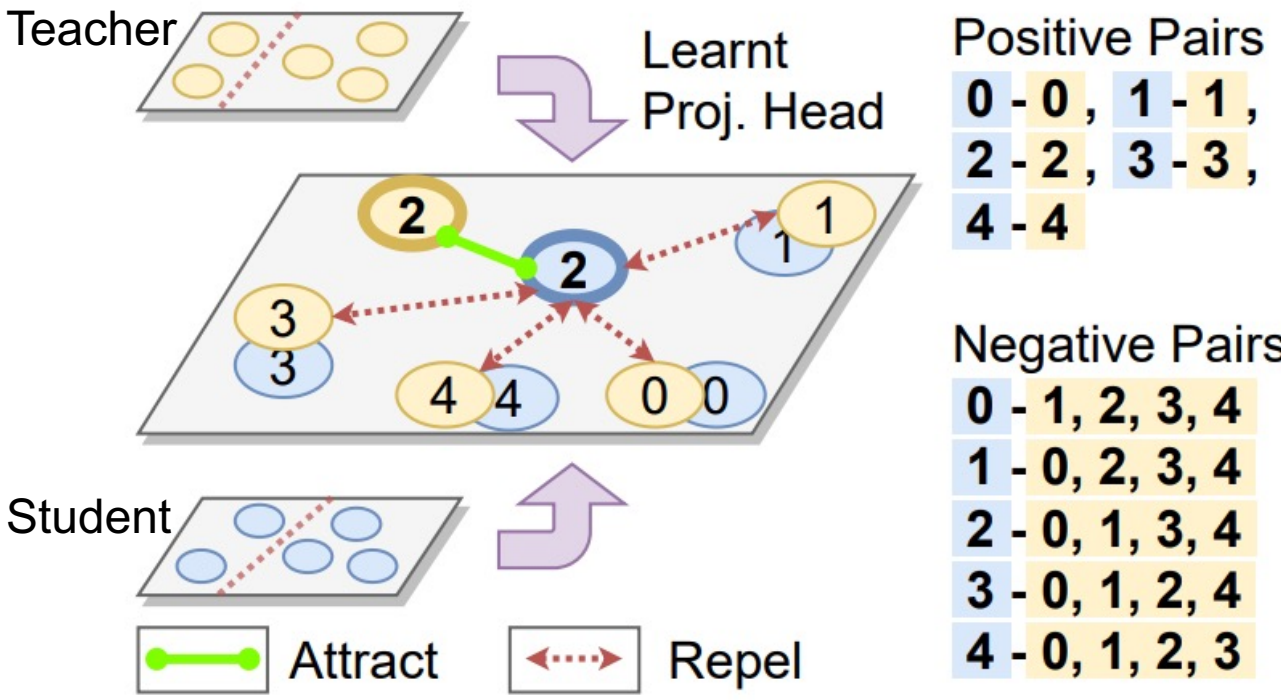


- **Global Structure Preserving (GSP):** consider all possible pairwise similarities among node features
- Challenging to optimize/scale up due to an explosion of possible pairs



# Knowledge Distillation from Structures

Enhance the student with implicit global structure: G-CRD



- Use contrastive learning to capture implicit global topology
- Positive pair: same node embeddings of the teacher and the student
- Negative pair: random nodes

# Knowledge Distillation from Structures

Enhance the student with implicit global structure: G-CRD



Molecular graph classification on MOLHIV

| Teacher (#Layer,#Param): |                    | GIN-E (5L,3.3M) | PNA (5L,2.4M)   | GIN-E (5L,3.3M) | PNA (5L,2.4M)   | PNA (5L,2.4M)   |
|--------------------------|--------------------|-----------------|-----------------|-----------------|-----------------|-----------------|
| Student (#Layer,#Param): |                    | GCN (2L,15K)    | GCN (2L,15K)    | GCN (2L,40K)    | GCN (2L,40K)    | GIN (2L,10K)    |
| Sup.                     | Supervised Teacher | 77.69 ±1.61     | 77.48 ±1.71     | 77.69 ±1.61     | 77.48 ±1.71     | 77.48 ±1.71     |
|                          | Supervised Student | 73.02 ±1.46     | 73.02 ±1.46     | 73.65 ±1.50     | 73.65 ±1.50     | 73.03 ±2.02     |
| Distillation             | KD [30]            | 74.08 ±1.03     | 74.13 ±1.72     | 75.25 ±1.71     | 74.45 ±1.27     | 73.42 ±2.14     |
|                          | FitNet [47]        | 73.62 ±1.05 (↓) | 73.65 ±1.25 (↓) | 74.52 ±1.33 (↓) | 74.39 ±1.46 (↓) | 72.88 ±0.89 (↓) |
|                          | AT [48]            | 73.85 ±0.85 (↓) | 73.64 ±1.50 (↓) | 74.94 ±0.97 (↓) | 73.89 ±1.92 (↓) | 73.87 ±2.28 (↑) |
|                          | LSP [32]           | 73.58 ±1.29 (↓) | 73.24 ±1.67 (↓) | 75.04 ±1.20 (↓) | 74.43 ±1.58 (↓) | 70.74 ±1.82 (↓) |
|                          | GSP                | 72.83 ±1.30 (↓) | 73.74 ±0.93 (↓) | 75.12 ±1.27 (↓) | 75.09 ±1.48 (↑) | 69.68 ±2.88 (↓) |
|                          | G-CRD (Ours)       | 74.34 ±1.44 (↑) | 75.11 ±0.73 (↑) | 75.53 ±1.64 (↑) | 75.89 ±0.80 (↑) | 75.77 ±2.02 (↑) |

- LSP performs poorly: preserving local structure is not sufficient
- GSP performs poorly: preserving every possible pair of nodes contributes less
- G-CRD performs well: preserving the implicit global structure is beneficial



# Covered Topics for Knowledge from Models



How to obtain knowledge from

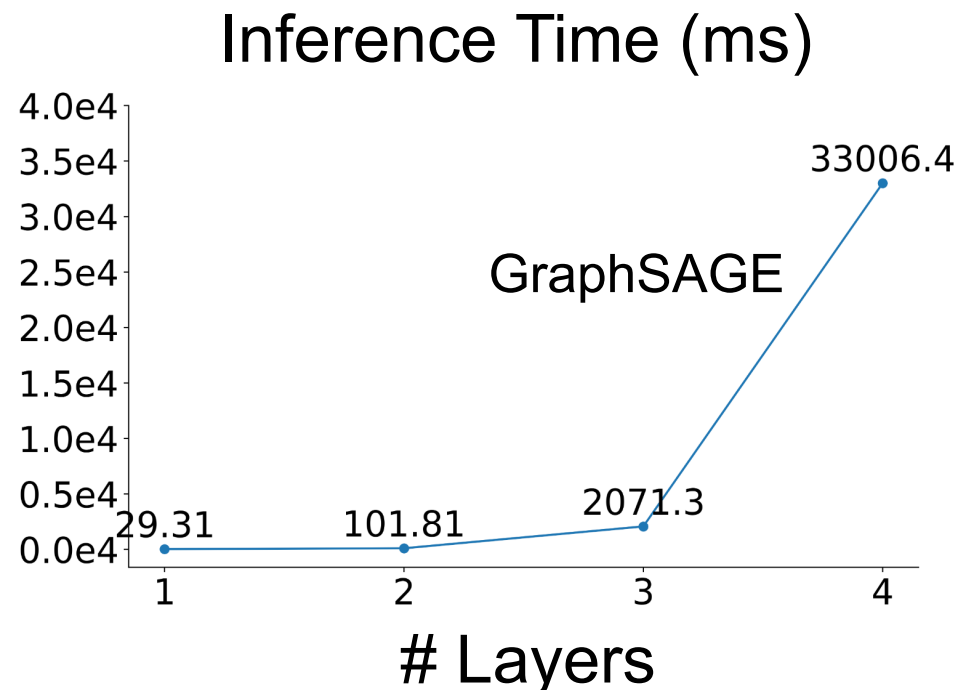
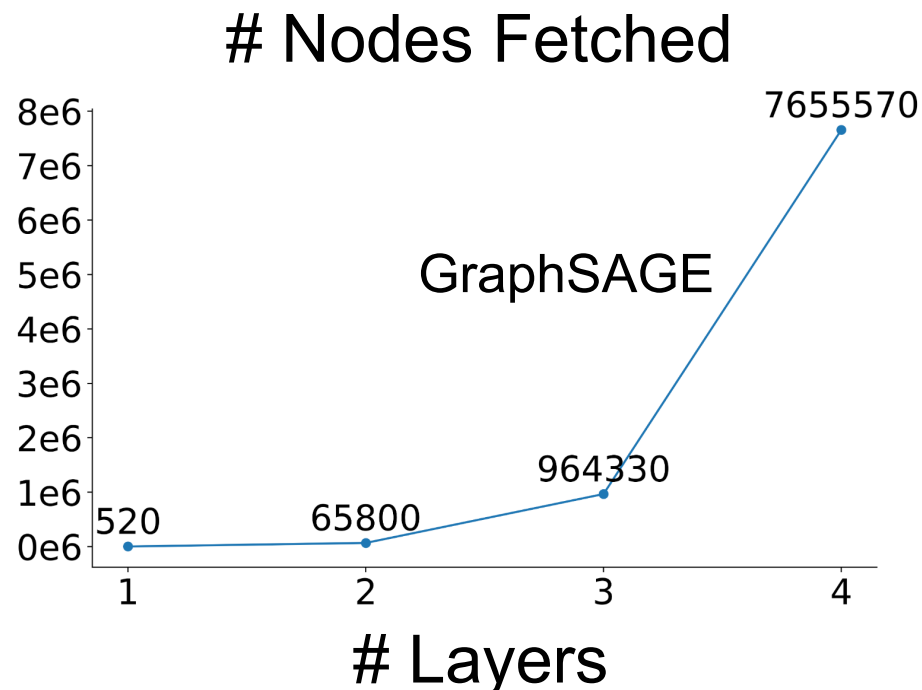
- 1) Logits
- 2) Embeddings
- 3) Structures

*A recent learning strategy*

 4) distilling an MLP to replace GNNs

# Distilling an MLP to Replace GNNs

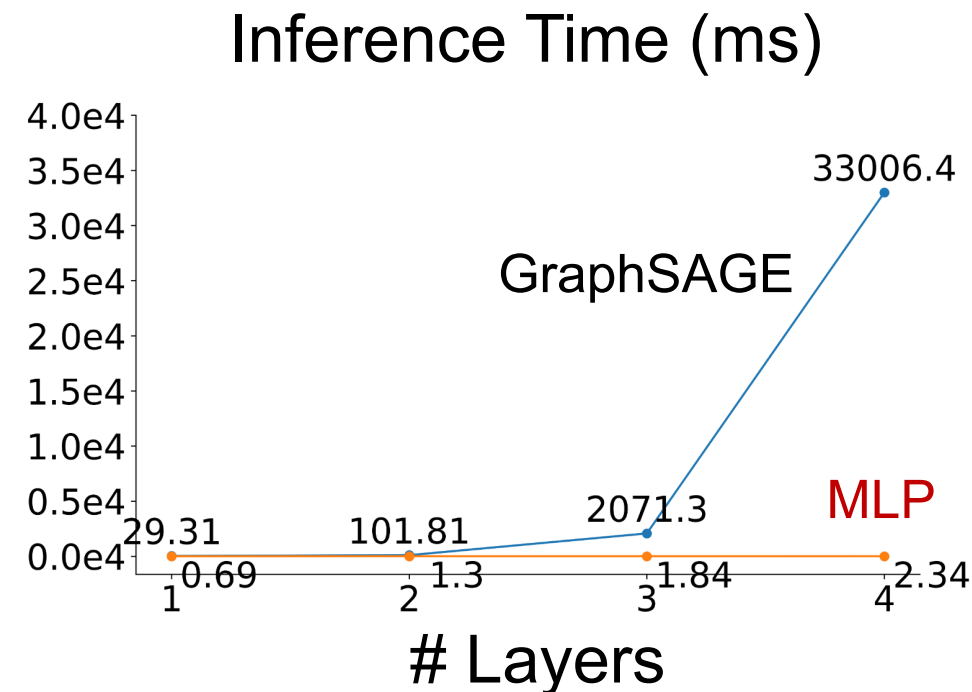
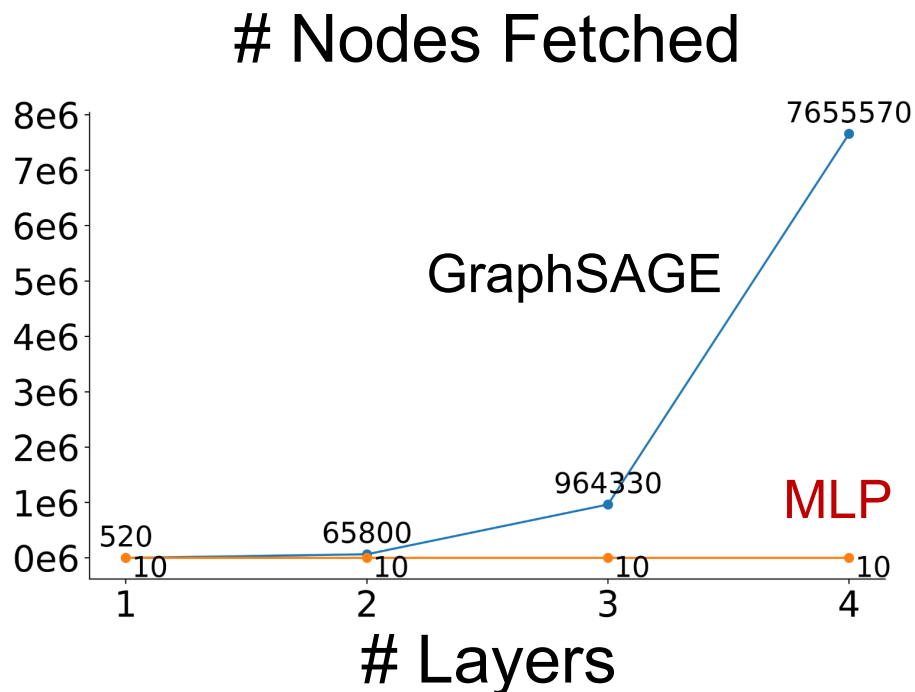
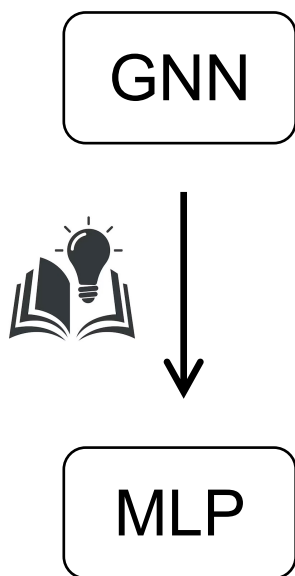
The majority of KD on Graphs focus on GNN to GNN



Inherent limitation: Dependence on message passing architecture, which is **time-consuming and computation-intensive**, making the model **inapplicable to time-sensitive situations**

# Distilling an MLP to Replace GNNs

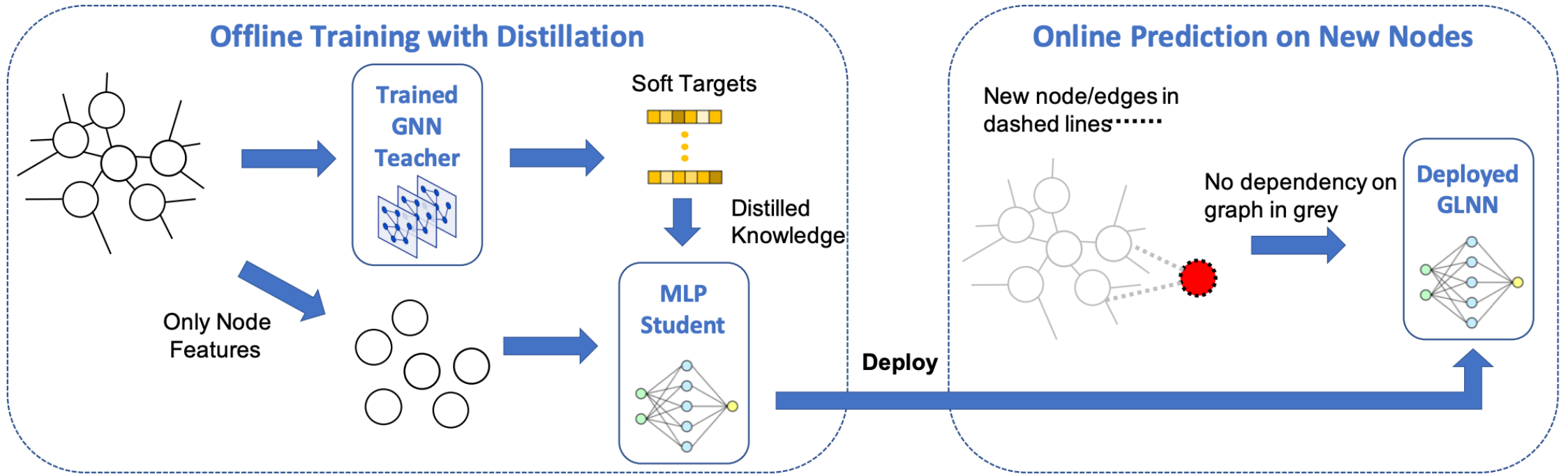
To avoid the troublesome message passing



**Distilling knowledge to an MLP** (Multi-layer Perceptron), which is simple and does not require message passing

# Distilling an MLP to Replace GNNs

## Distilling knowledge to an MLP: GLNN



# Distilling an MLP to Replace GNNs

## Distilling knowledge to an MLP: GLNN



| Datasets   | SAGE                               | MLP              | GLNN                               | $\Delta_{MLP}$ | $\Delta_{GNN}$  |
|------------|------------------------------------|------------------|------------------------------------|----------------|-----------------|
| Cora       | 80.52 $\pm$ 1.77                   | 59.22 $\pm$ 1.31 | <b>80.54 <math>\pm</math> 1.35</b> | 21.32 (36.00%) | 0.02 (0.02%)    |
| Citeseer   | 70.33 $\pm$ 1.97                   | 59.61 $\pm$ 2.88 | <b>71.77 <math>\pm</math> 2.01</b> | 12.16 (20.40%) | 1.44 (2.05%)    |
| Pubmed     | 75.39 $\pm$ 2.09                   | 67.55 $\pm$ 2.31 | <b>75.42 <math>\pm</math> 2.31</b> | 7.87 (11.65%)  | 0.03 (0.04%)    |
| A-computer | 82.97 $\pm$ 2.16                   | 67.80 $\pm$ 1.06 | <b>83.03 <math>\pm</math> 1.87</b> | 15.23 (22.46%) | 0.06 (0.07%)    |
| A-photo    | 90.90 $\pm$ 0.84                   | 78.77 $\pm$ 1.74 | <b>92.11 <math>\pm</math> 1.08</b> | 13.34 (16.94%) | 1.21 (1.33%)    |
| Arxiv      | <b>70.92 <math>\pm</math> 0.17</b> | 56.05 $\pm$ 0.46 | 63.46 $\pm$ 0.45                   | 7.41 (13.24%)  | -7.46 (-10.52%) |
| Products   | <b>78.61 <math>\pm</math> 0.49</b> | 62.47 $\pm$ 0.10 | 68.86 $\pm$ 0.46                   | 6.39 (10.23%)  | -9.75 (-12.4%)  |

GLNN performs well on small datasets,  
but performs poorly on large datasets

# Distilling an MLP to Replace GNNs



## Problems of simply distilling knowledge to an MLP

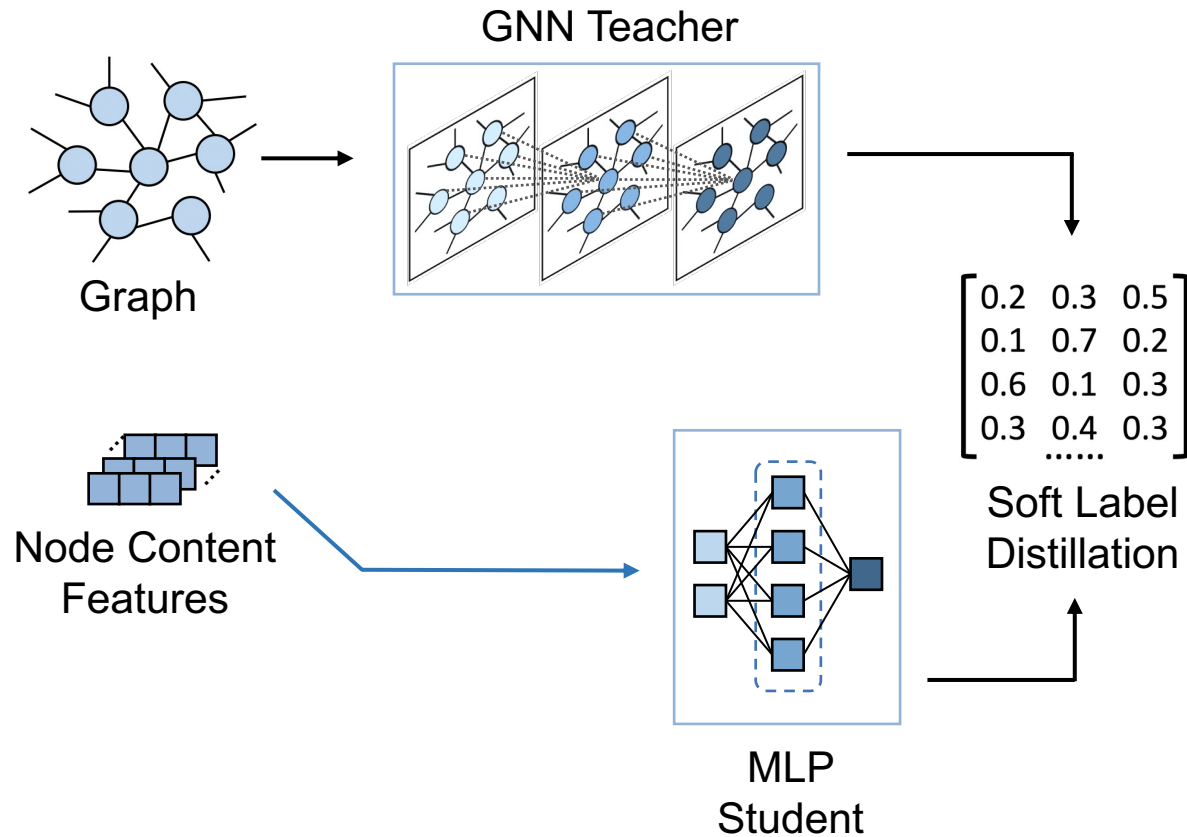
- The misalignment between content feature and label spaces
- The strict hard matching to teacher's output
- The sensitivity to node feature noises

**Therefore, a question can be asked:**

Can we learn MLPs that are graph structure-aware in both the feature and representation spaces, insensitive to node feature noises, and have superior performance as well as fast inference speed?

# Distilling an MLP to Replace GNNs

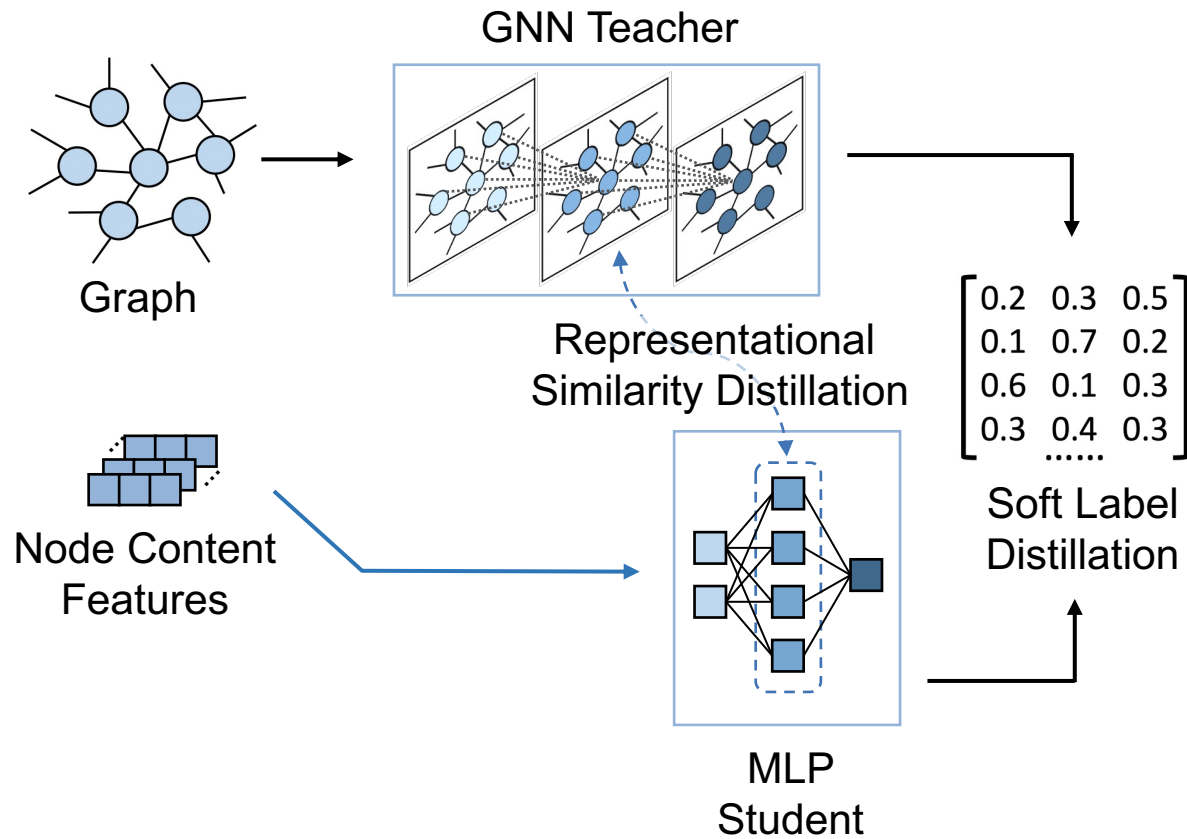
Learn a Noise-robust Structure-aware MLPs On Graphs: NOSMOG



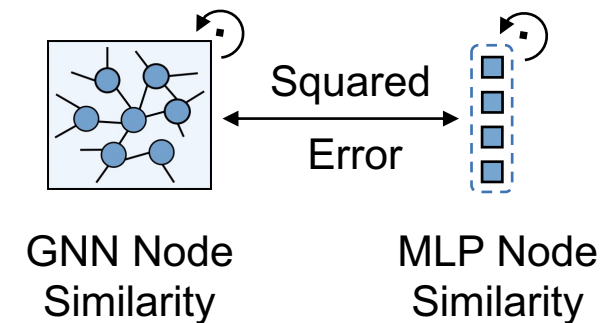
- Distill Knowledge of **Logits**

# Distilling an MLP to Replace GNNs

Learn a Noise-robust Structure-aware MLPs On Graphs: NOSMOG



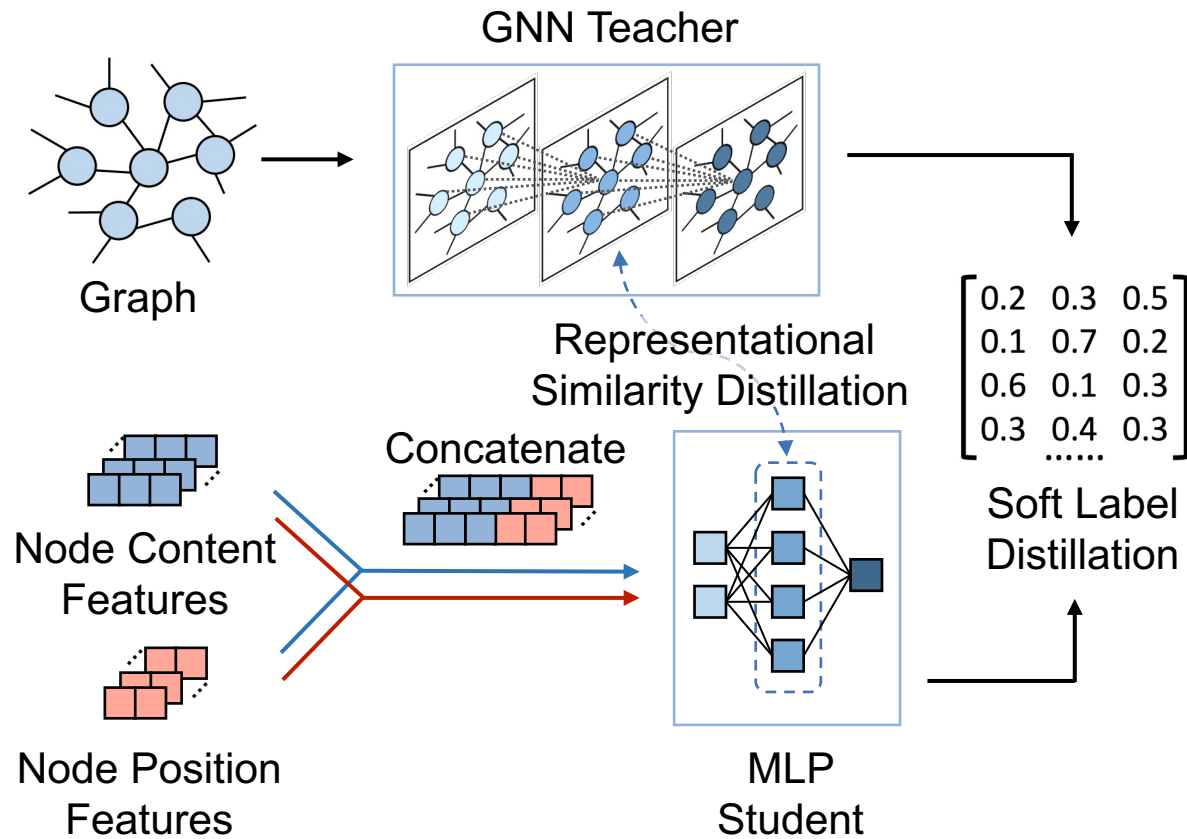
- Distill Knowledge of Logits
- Distill Knowledge of **Learned Embedding and Structure**





# Distilling an MLP to Replace GNNs

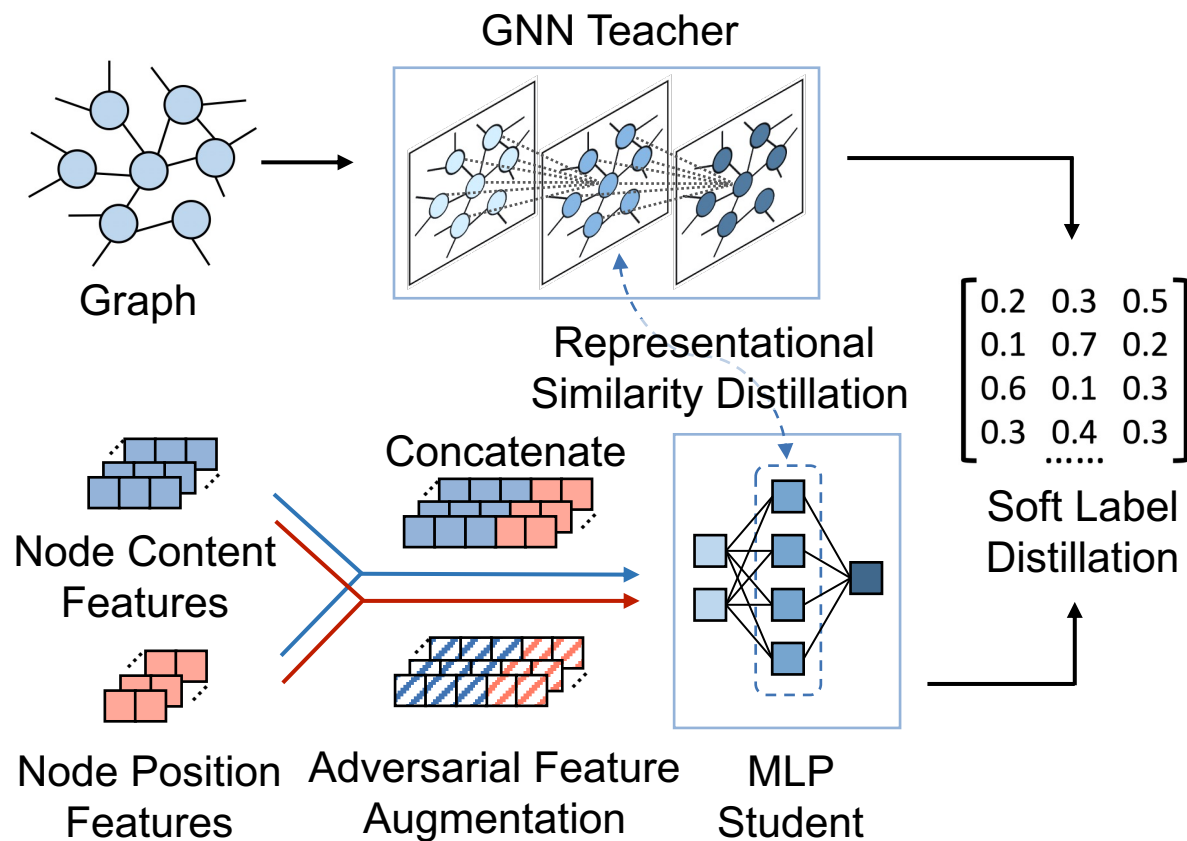
Learn a Noise-robust Structure-aware MLPs On Graphs: NOSMOG



- Distill Knowledge of Logits
- Distill Knowledge of Learned Embedding and Structure
- Learning Knowledge of **Node Positions** from Data

# Distilling an MLP to Replace GNNs

Learn a Noise-robust Structure-aware MLPs On Graphs: NOSMOG



- Distill Knowledge of Logits
- Distill Knowledge of Learned Embedding and Structure
- Learning Knowledge of Node Positions from Data
- **Enhance Knowledge Acquisition with Augmentation**

# Distilling an MLP to Replace GNNs

Learn a Noise-robust Structure-aware MLPs On Graphs: NOSMOG



| Datasets   | SAGE             | MLP              | GLNN             | NOSMOG                             | $\Delta_{GNN}$   | $\Delta_{MLP}$    | $\Delta_{GLNN}$   |
|------------|------------------|------------------|------------------|------------------------------------|------------------|-------------------|-------------------|
| Cora       | 80.64 $\pm$ 1.57 | 59.18 $\pm$ 1.60 | 80.26 $\pm$ 1.66 | <b>83.04 <math>\pm</math> 1.26</b> | $\uparrow$ 2.98% | $\uparrow$ 40.32% | $\uparrow$ 3.46%  |
| Citeseer   | 70.49 $\pm$ 1.53 | 58.50 $\pm$ 1.86 | 71.22 $\pm$ 1.50 | <b>73.78 <math>\pm</math> 1.54</b> | $\uparrow$ 4.67% | $\uparrow$ 26.12% | $\uparrow$ 3.59%  |
| Pubmed     | 75.56 $\pm$ 2.06 | 68.39 $\pm$ 3.09 | 75.59 $\pm$ 2.46 | <b>77.34 <math>\pm</math> 2.36</b> | $\uparrow$ 2.36% | $\uparrow$ 13.09% | $\uparrow$ 2.32%  |
| A-computer | 82.82 $\pm$ 1.37 | 67.62 $\pm$ 2.21 | 82.71 $\pm$ 1.18 | <b>84.04 <math>\pm</math> 1.01</b> | $\uparrow$ 1.47% | $\uparrow$ 24.28% | $\uparrow$ 1.61%  |
| A-photo    | 90.85 $\pm$ 0.87 | 77.29 $\pm$ 1.79 | 91.95 $\pm$ 1.04 | <b>93.36 <math>\pm</math> 0.69</b> | $\uparrow$ 2.76% | $\uparrow$ 20.79% | $\uparrow$ 1.53%  |
| Arxiv      | 70.73 $\pm$ 0.35 | 55.67 $\pm$ 0.24 | 63.75 $\pm$ 0.48 | <b>71.65 <math>\pm</math> 0.29</b> | $\uparrow$ 1.30% | $\uparrow$ 28.70% | $\uparrow$ 12.39% |
| Products   | 77.17 $\pm$ 0.32 | 60.02 $\pm$ 0.10 | 63.71 $\pm$ 0.31 | <b>78.45 <math>\pm</math> 0.38</b> | $\uparrow$ 1.66% | $\uparrow$ 30.71% | $\uparrow$ 23.14% |

+2.46%    +26.29%    +6.86%

NOSMOG achieves the best performance

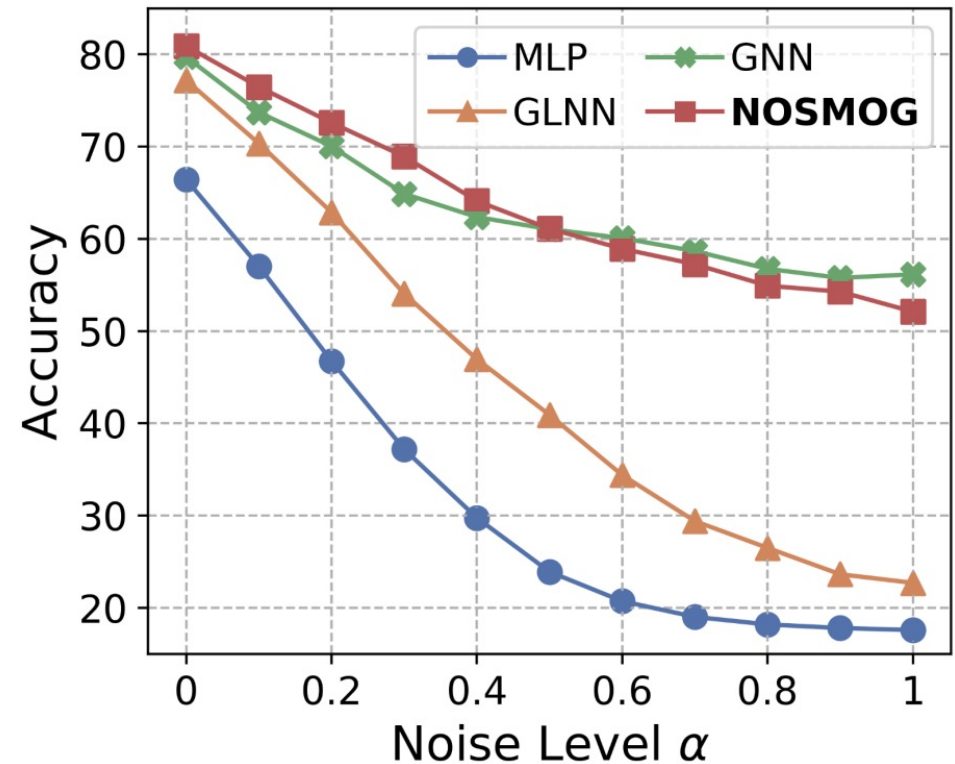
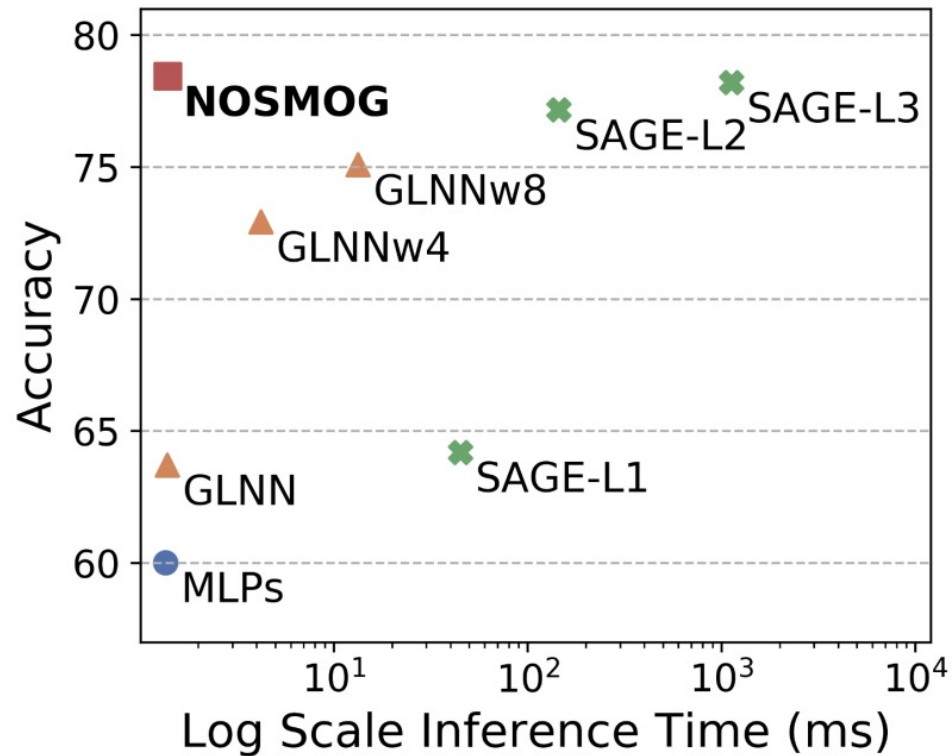
# Distilling an MLP to Replace GNNs

Learn a Noise-robust Structure-aware MLPs On Graphs: NOSMOG



NOSMOG is 833x faster than GNNs

NOSMOG is as robust as GNNs



# Graph Learning Enhanced by Knowledge from Models



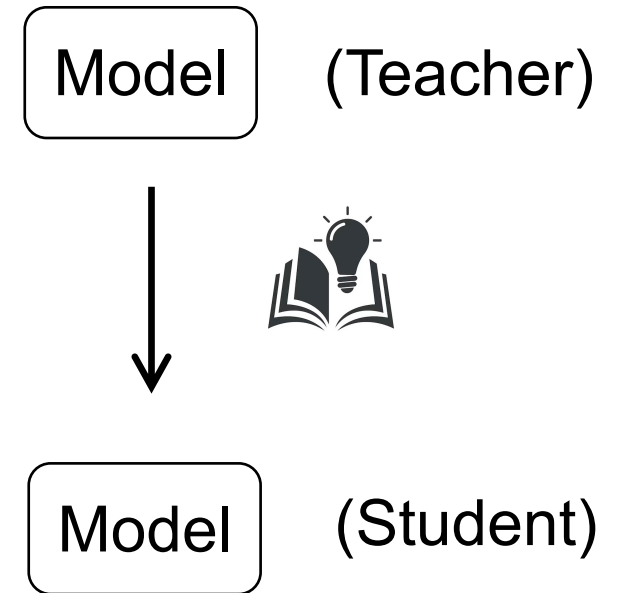
## Chapter summary

How to obtain knowledge from

- 1) Logits
- 2) Embeddings
- 3) Structures

A recent learning strategy

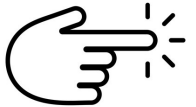
- 4) distilling an MLP to replace GNNs



# Tutorial Outline



- Preliminaries and Foundations
- Graph Learning Enhanced by Knowledge from Data
- Graph Learning Enhanced by Knowledge from Models
- **Graph Learning Enhanced by Knowledge from Humans and Domains**
- Graph Learning Enhanced by Knowledge from External Sources
- Knowledge-enhanced Graph Learning for Real-world Applications
- Summary and Future Directions



# Topics for Knowledge from Humans and Domains



## Graph learning enhanced by human feedback

- Graph active learning with human knowledge

## Graph learning enhanced by domain knowledge

- Chemistry domain knowledge for molecular property prediction



# Topics for Knowledge from Humans and Domains



Graph learning enhanced by human feedback



- Graph active learning with human knowledge

Graph learning enhanced by domain knowledge

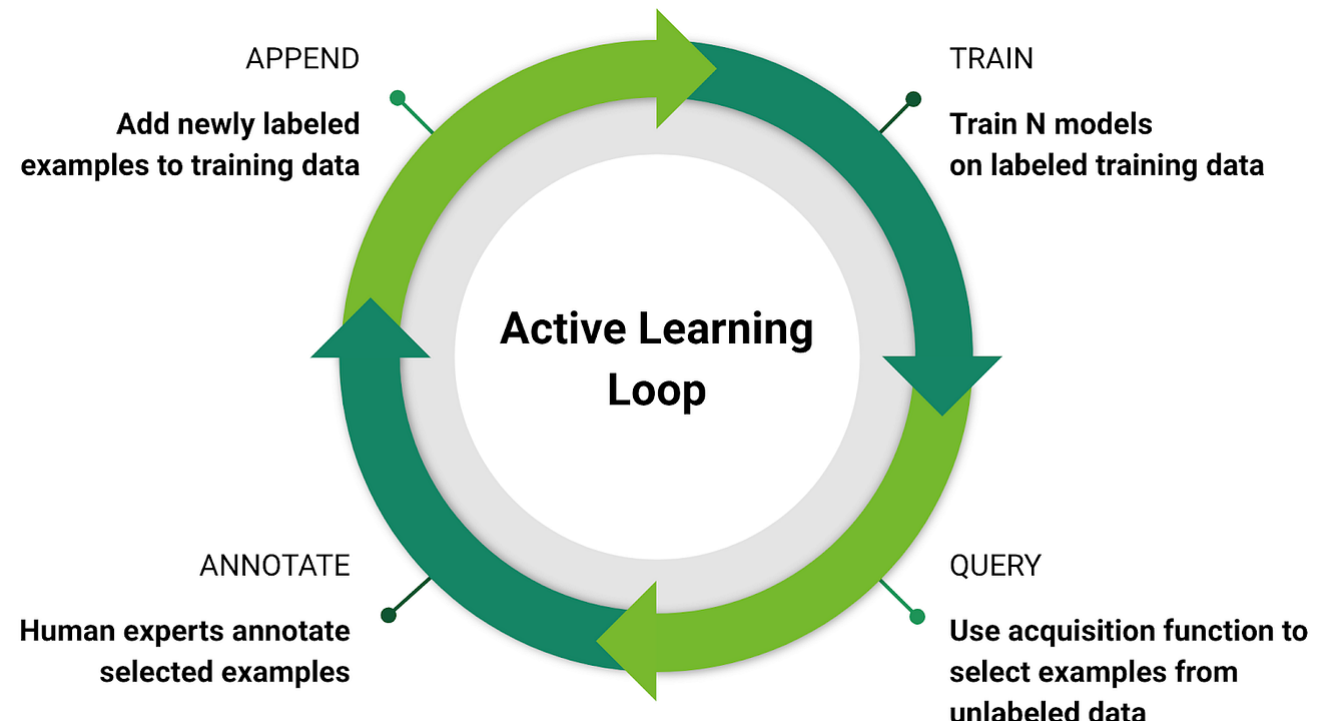
- Chemistry domain knowledge for molecular property prediction

# Graph Active Learning

## What is graph active learning?

Chooses graph data examples to label from a large pool of unlabeled data repeatedly:

- training a model on the small pool of labeled graph data
- selecting graph data examples to label based on different query heuristics

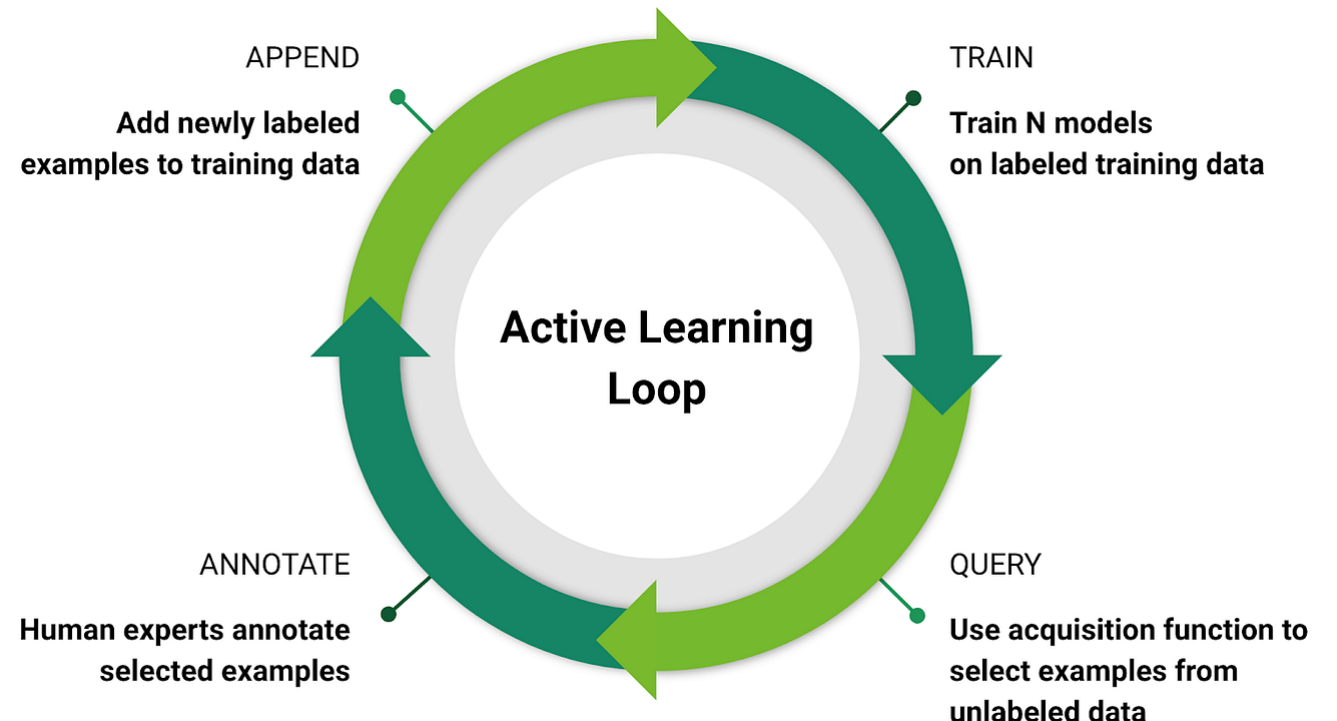


# Graph Active Learning

## Why do we need graph active learning?

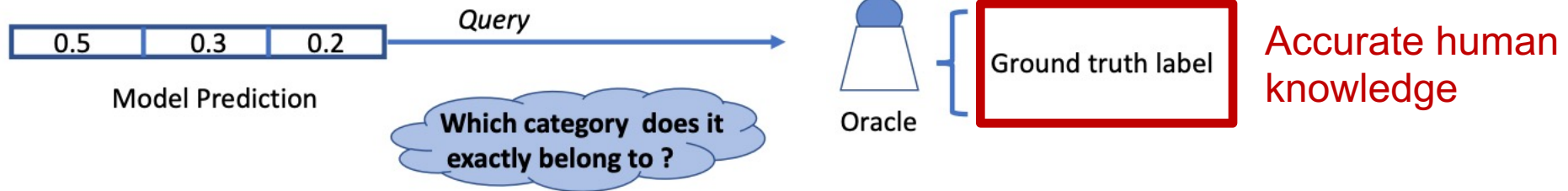


- Incorporating human knowledge
- Reduce the annotation cost by focusing on the most relevant graph data examples



# Graph Active Learning

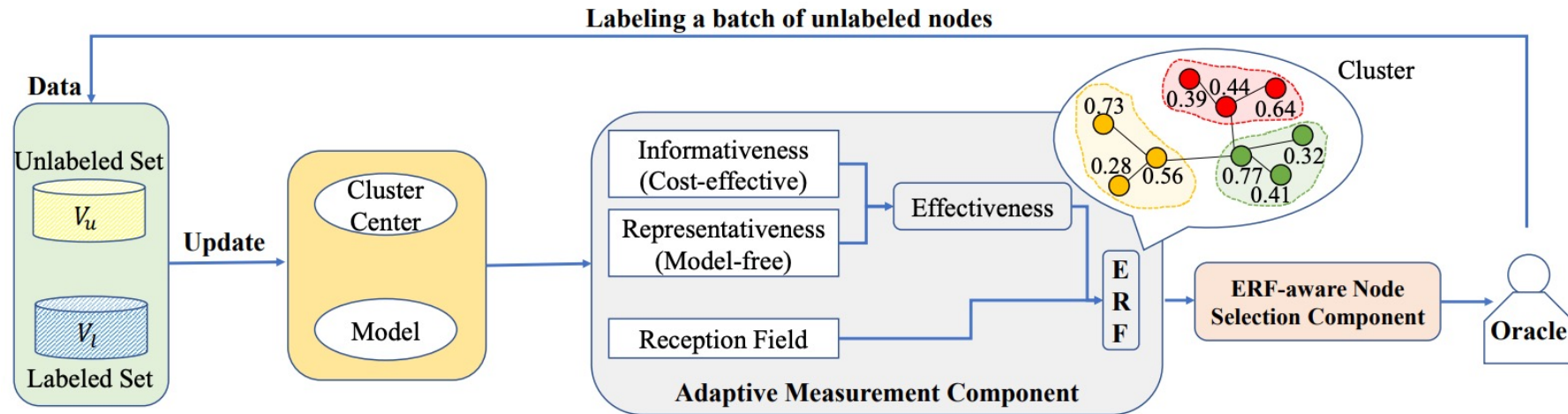
## Incorporating accurate human knowledge: ALG



In this setting, ground truth labels that are one-hot encoding serve as accurate human knowledge

# Graph Active Learning

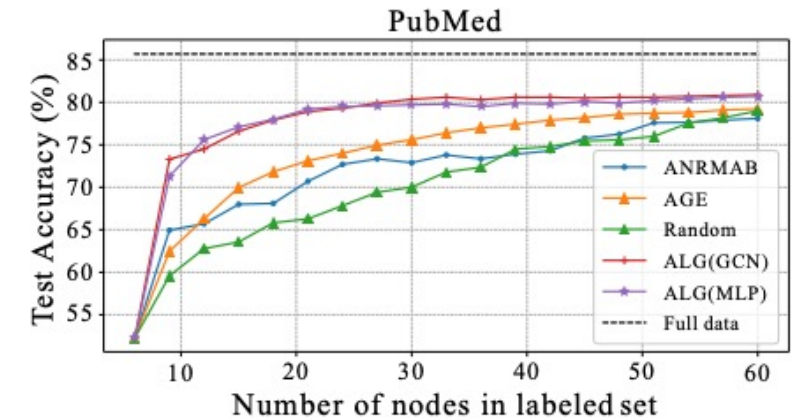
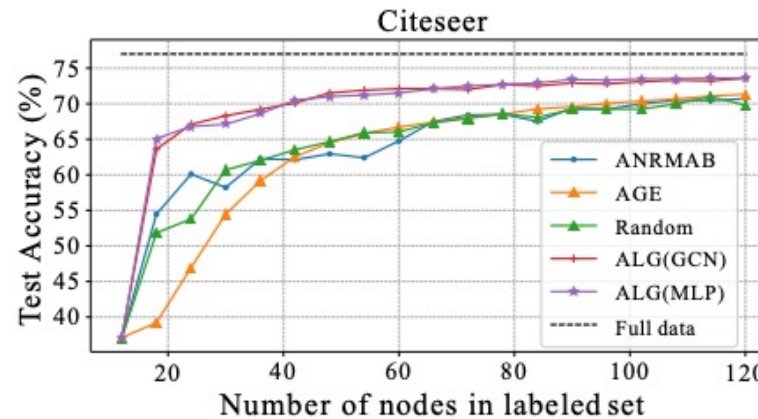
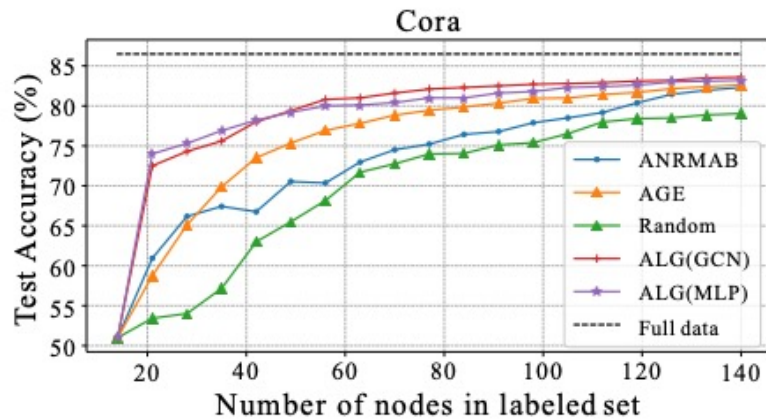
## Incorporating accurate human knowledge: ALG



- An adaptive measurement component that computes the node importance, and a selection component to choose a set of nodes to label given their importance
- Measurement component: K-means-based. Center nodes in clusters are more representative
- Selection component: the model should select nodes from each cluster for fairness

# Graph Active Learning

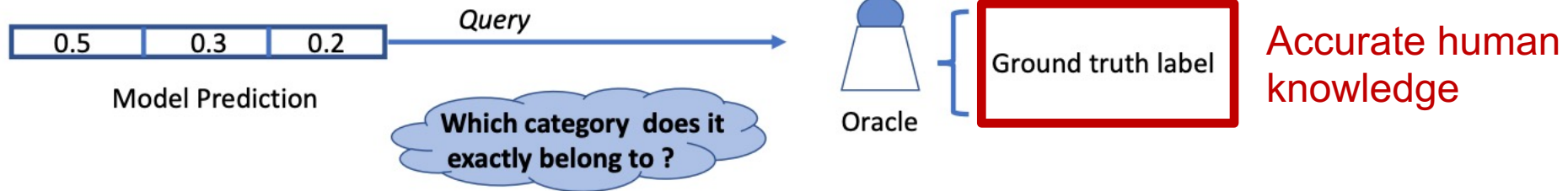
## Incorporating accurate human knowledge: ALG



- When labeling more nodes, ALG quickly boosts its accuracy at the beginning and consistently outperforms the baselines
- ALG only needs to label about 5% of all the nodes in Cora to achieve the accuracy of 83.5%, which is comparable to the performance trained on the full dataset (a gap <3%)

# Graph Active Learning

## Problem of using accurate human knowledge

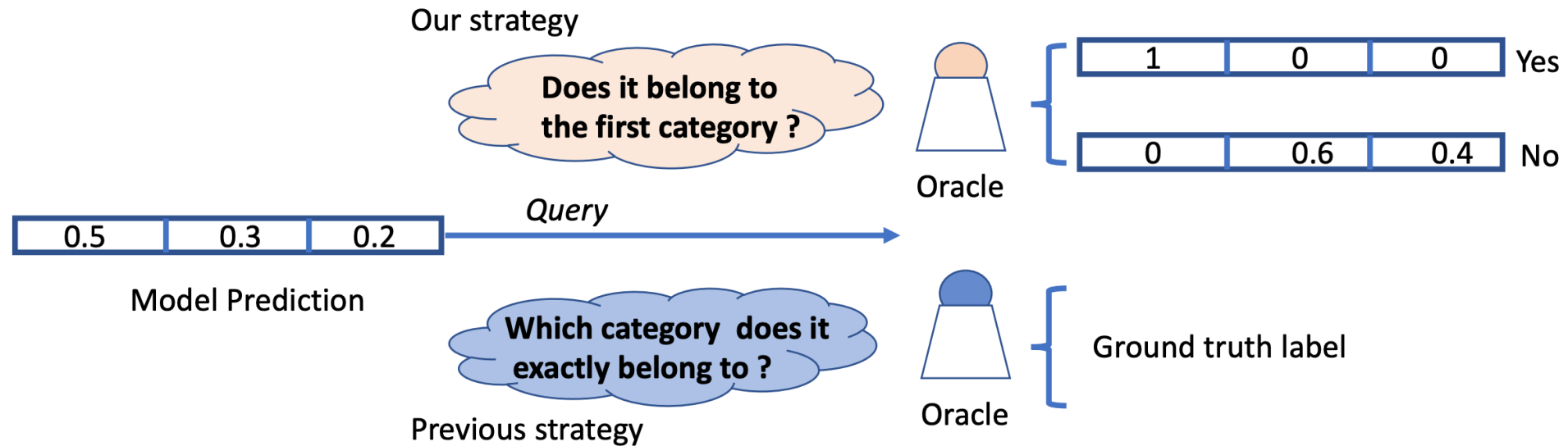


Exact labeling task (for ground truth labels) is costly, especially when the categorization task is specialized and strongly depend on the oracle expertise



# Graph Active Learning

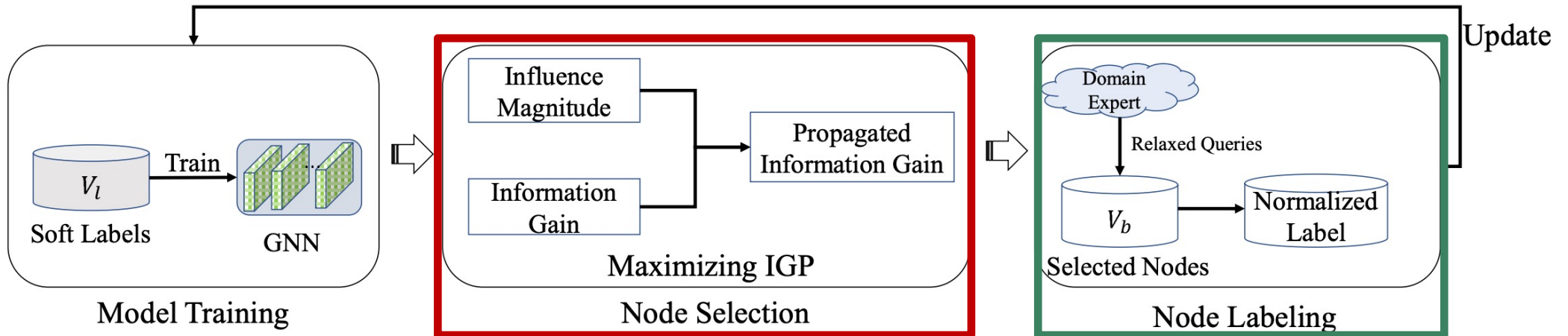
## Human knowledge as soft label: IGP



- A domain expert only judges the correctness of the predicted labels (a binary question) rather than identifying the exact class (a multi-class question)
- If the model prediction is incorrect, the node is annotated with the soft label by re-normalizing the model outputs over the remaining classes

# Graph Active Learning

## Human knowledge as soft label: IGP



IGP measures the **expected information gain** of labeling each node and selects a batch of nodes that can maximize the information gain propagation on the graph

The oracle only needs to judge the **correctness of the generated label**

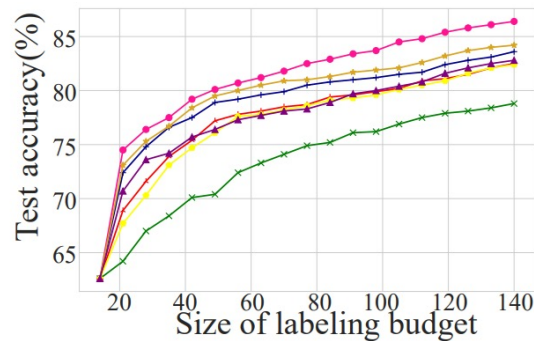
# Graph Active Learning

## Human knowledge as soft label: IGP

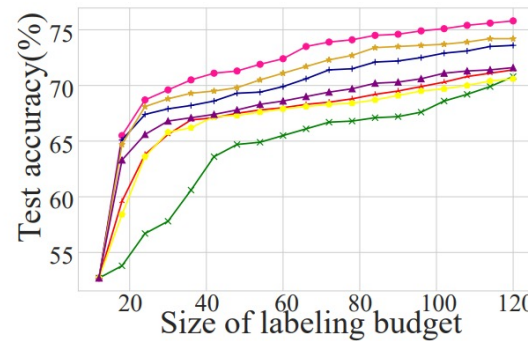


| Method     | Cora                              | Citeseer                          | PubMed                            | Reddit                            | ogbn-arxiv                        |
|------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|
| Random     | 78.8( $\pm 0.8$ )                 | 70.8( $\pm 0.9$ )                 | 78.9( $\pm 0.6$ )                 | 91.1( $\pm 0.5$ )                 | 68.2( $\pm 0.4$ )                 |
| AGE        | 82.5( $\pm 0.6$ )                 | 71.4( $\pm 0.6$ )                 | 79.4( $\pm 0.4$ )                 | 91.6( $\pm 0.3$ )                 | 68.9( $\pm 0.3$ )                 |
| ANRMAB     | 82.4( $\pm 0.5$ )                 | 70.6( $\pm 0.6$ )                 | 78.2( $\pm 0.3$ )                 | 91.5( $\pm 0.3$ )                 | 68.7( $\pm 0.2$ )                 |
| GPA        | 82.8( $\pm 0.4$ )                 | 71.6( $\pm 0.4$ )                 | 79.9( $\pm 0.5$ )                 | 91.8( $\pm 0.2$ )                 | 69.2( $\pm 0.3$ )                 |
| SEAL       | 83.2( $\pm 0.5$ )                 | 72.1( $\pm 0.4$ )                 | 80.3( $\pm 0.4$ )                 | 92.1( $\pm 0.3$ )                 | 69.5( $\pm 0.1$ )                 |
| ALG        | 83.6( $\pm 0.6$ )                 | 73.6( $\pm 0.5$ )                 | 80.9( $\pm 0.3$ )                 | 92.4( $\pm 0.3$ )                 | 70.1( $\pm 0.2$ )                 |
| GRAIN      | 84.2( $\pm 0.3$ )                 | 74.2( $\pm 0.3$ )                 | 81.8( $\pm 0.2$ )                 | 92.5( $\pm 0.1$ )                 | 70.3( $\pm 0.2$ )                 |
| <b>IGP</b> | <b>86.4(<math>\pm 0.6</math>)</b> | <b>75.8(<math>\pm 0.3</math>)</b> | <b>83.6(<math>\pm 0.5</math>)</b> | <b>93.4(<math>\pm 0.2</math>)</b> | <b>70.9(<math>\pm 0.3</math>)</b> |

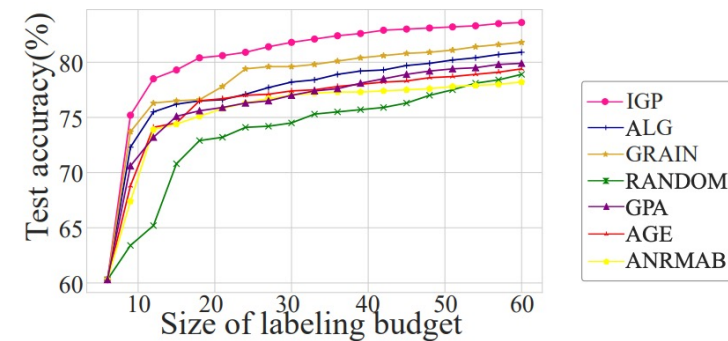
Better performance on node classification



(a) Cora

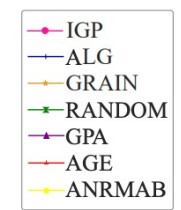


(b) Citeseer



(c) PubMed

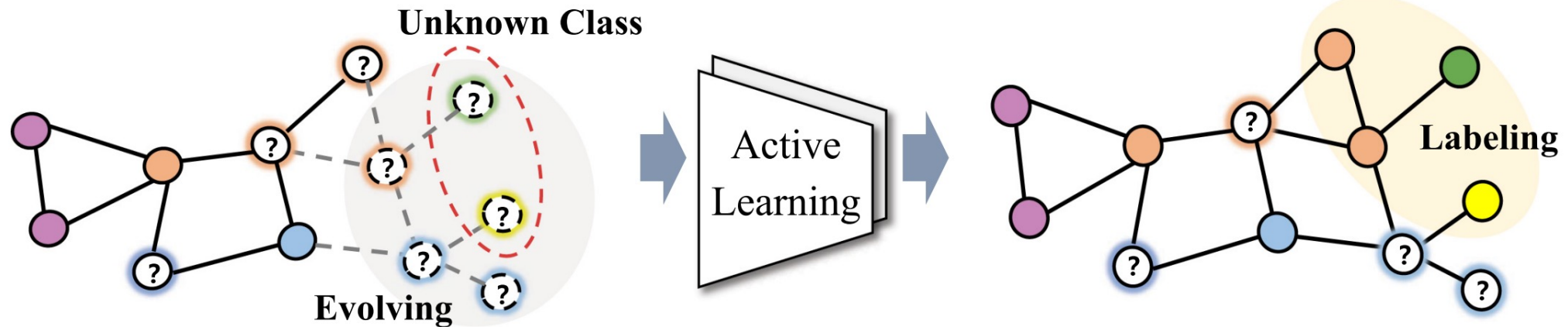
Different labeling budgets



# Graph Active Learning

## Open world scenario: OWGAL

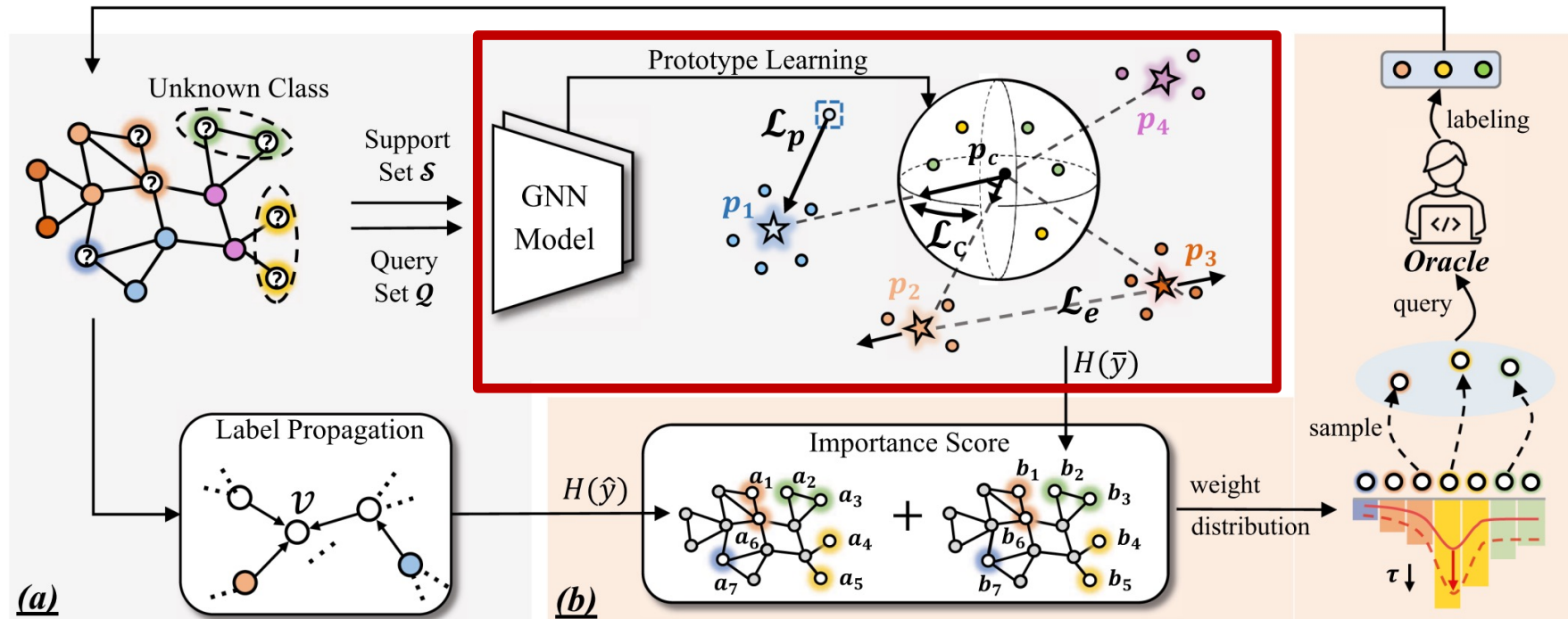
- Existing GAL methods mainly focus on a “closed-world” setting, where all nodes belong to a fixed known group of classes
- In the open world scenario, a graph develops over time as new nodes, edges, and classes are introduced, and thus GAL is required to make a fair selection among all nodes for labeling, despite its class has been seen or not



Open-world graph active learning (OWGAL) should not only select the most informative nodes, but uncover nodes of latent new classes

# Graph Active Learning

## Open world scenario: OWGAL

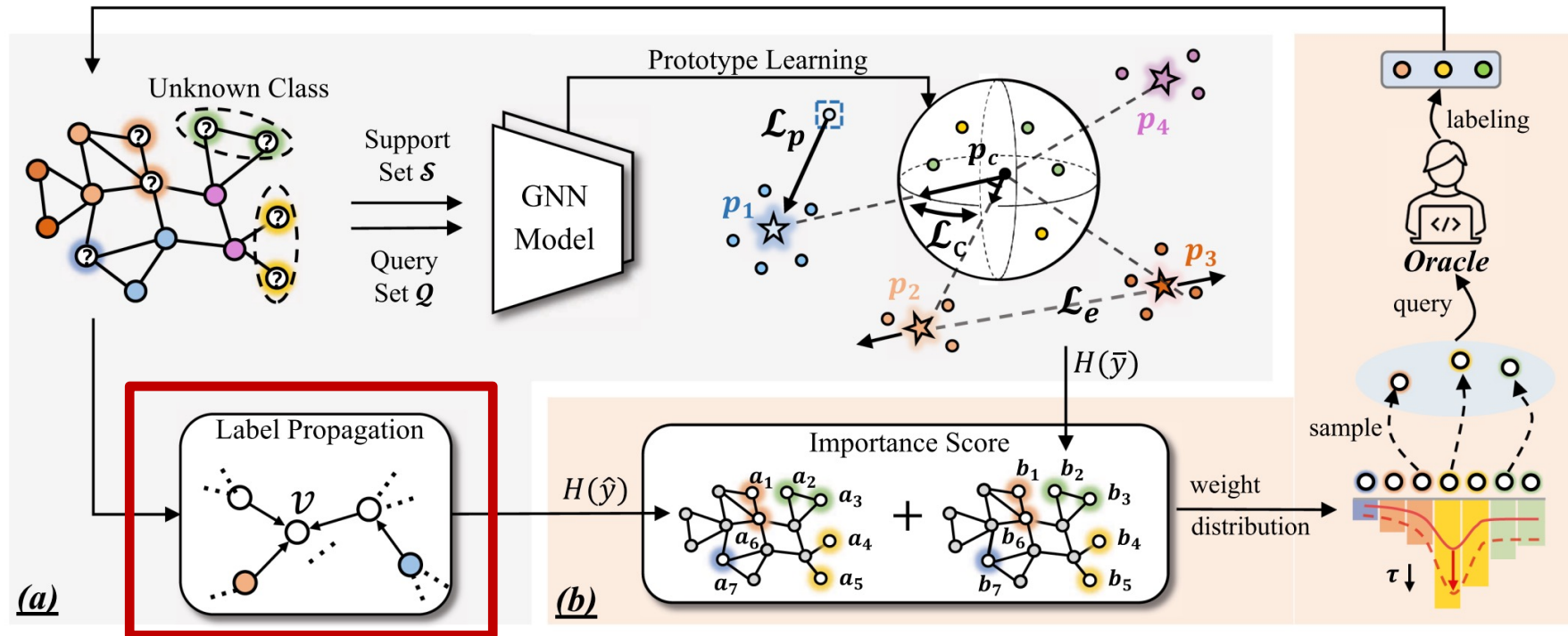


To make known and unknown classes more identifiable, learn more compact representations for nodes in the same class, and push known classes as far away from each other as possible



# Graph Active Learning

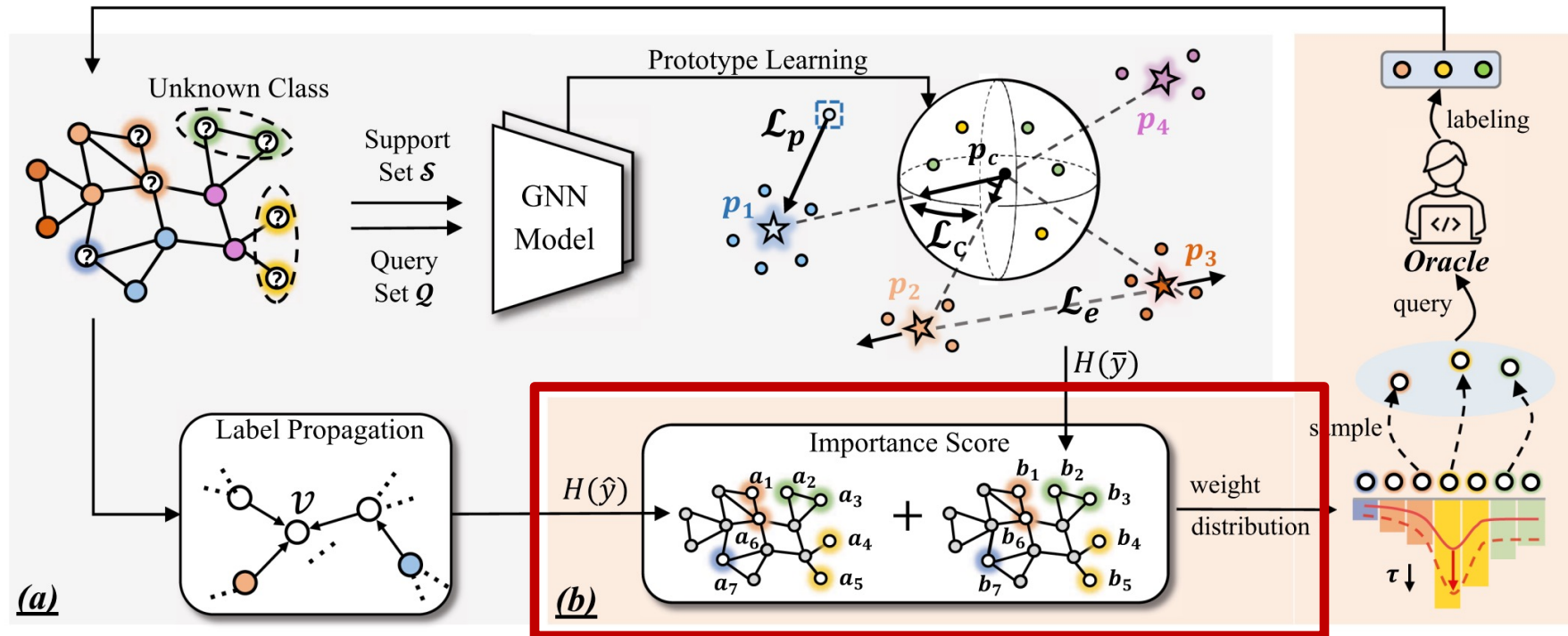
## Open world scenario: OWGAL



Label propagation: motivated by the homophily assumption that adjacent nodes tend to have the same label, node  $v$  is more likely to belong to a novel class if no labeled node reach it within  $K$  steps of random walks on graph

# Graph Active Learning

## Open world scenario: OWGAL



Unify the discovery of unknown classes and important nodes, as it weighs unseen-class nodes as much as the nodes near the decision boundary, both of which have high uncertainty scores and should be sampled for labeling



# Graph Active Learning

## Open world scenario: OWGAL



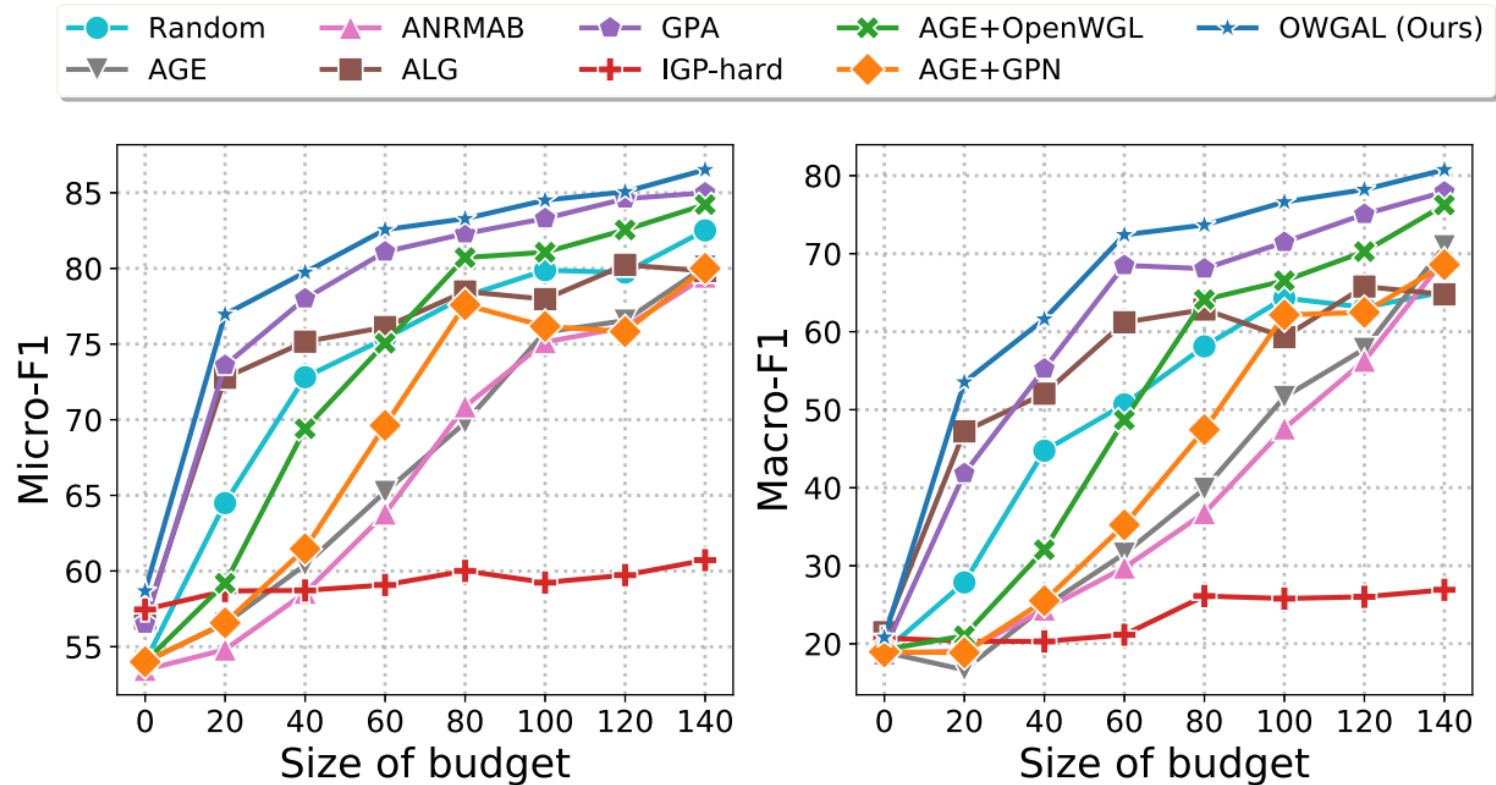
| Baselines           | Amazon-Computer      |                      | Amazon-Photo         |                      | Coauthor-CS         |                     | Coauthor-Physics     |                      |
|---------------------|----------------------|----------------------|----------------------|----------------------|---------------------|---------------------|----------------------|----------------------|
|                     | Micro-F1             | Macro-F1             | Micro-F1             | Macro-F1             | Micro-F1            | Macro-F1            | Micro-F1             | Macro-F1             |
| Random              | 82.51 ± 1.77%        | 74.16 ± 2.90%        | 89.88 ± 1.55%        | 86.65 ± 2.94%        | 88.25 ± 0.34%       | 80.18 ± 1.03%       | 91.43 ± 1.01%        | 87.50 ± 1.83%        |
| AGE                 | 80.12 ± 2.24%        | 70.94 ± 3.08%        | 90.35 ± 0.87%        | 89.17 ± 1.32%        | 89.08 ± 1.69%       | 82.69 ± 2.05%       | 91.76 ± 2.30%        | 88.08 ± 3.99%        |
| ANRMAB              | 79.57 ± 1.53%        | 69.28 ± 2.77%        | 89.63 ± 0.53%        | 88.55 ± 1.07%        | 89.87 ± 1.32%       | 82.33 ± 1.55%       | 90.96 ± 2.02%        | 87.68 ± 1.82%        |
| GPA                 | 82.92 ± 1.45%        | 75.31 ± 3.71%        | <u>91.59 ± 0.41%</u> | <u>90.75 ± 0.45%</u> | 88.89 ± 1.51%       | 85.78 ± 1.82%       | 91.08 ± 1.49%        | 87.86 ± 2.22%        |
| ALG                 | 79.83 ± 1.85%        | 64.79 ± 5.91%        | 80.61 ± 1.30%        | 83.58 ± 2.64%        | 85.97 ± 0.09%       | 72.01 ± 0.19%       | 90.86 ± 0.28%        | 87.96 ± 0.74%        |
| IGP-hard            | 65.72 ± 2.37%        | 31.92 ± 1.55%        | 81.47 ± 1.61%        | 80.23 ± 1.71%        | 82.34 ± 0.23%       | 59.51 ± 0.43%       | 85.92 ± 1.42%        | 65.28 ± 3.66%        |
| AGE+OpenWGL         | <u>84.21 ± 0.27%</u> | <u>76.24 ± 1.05%</u> | 91.33 ± 1.14%        | 90.00 ± 1.36%        | 88.03 ± 0.19%       | 76.26 ± 0.51%       | 92.79 ± 0.34%        | 89.74 ± 0.54%        |
| AGE+GPN             | 80.01 ± 2.22%        | 68.57 ± 4.10%        | 89.93 ± 0.52%        | 88.98 ± 0.62%        | 90.16 ± 0.09%       | 87.03 ± 0.25%       | <u>92.92 ± 0.19%</u> | <u>90.11 ± 0.23%</u> |
| <b>OWGAL (Ours)</b> | <b>86.39 ± 1.20%</b> | <b>80.79 ± 2.80%</b> | <b>93.09 ± 0.61%</b> | <b>91.71 ± 0.79%</b> | <b>91.42 ± 0.4%</b> | <b>88.77 ± 0.5%</b> | <b>93.11 ± 0.28%</b> | <b>90.36 ± 0.62%</b> |

| Baselines           | Arxiv                |                      | Reddit               |                      | Products             |                      |
|---------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|
|                     | Micro-F1             | Macro-F1             | Micro-F1             | Macro-F1             | Micro-F1             | Macro-F1             |
| Random              | <u>62.53 ± 0.74%</u> | 33.79 ± 1.12%        | 85.54 ± 0.94%        | <u>70.47 ± 2.25%</u> | <u>63.61 ± 0.68%</u> | 22.89 ± 0.59%        |
| AGE                 | 61.44 ± 0.38%        | 31.91 ± 1.44%        | 78.40 ± 0.43%        | 55.39 ± 1.09%        | 62.48 ± 1.04%        | <u>23.29 ± 0.69%</u> |
| ANRMAB              | 61.85 ± 0.54%        | 30.67 ± 1.83%        | 77.40 ± 0.58%        | 56.12 ± 1.36%        | 62.55 ± 0.93%        | 23.07 ± 0.83%        |
| GPA                 | 58.59 ± 0.34%        | 31.34 ± 0.98%        | 85.00 ± 1.52%        | 69.65 ± 5.96%        | 58.35 ± 1.06%        | 21.22 ± 0.25%        |
| ALG                 | 59.71 ± 0.50%        | 22.36 ± 1.28%        | <u>86.14 ± 0.67%</u> | 68.55 ± 1.27%        | 63.59 ± 0.88%        | 21.63 ± 0.67%        |
| IGP-hard            | 56.35 ± 3.58%        | 22.92 ± 2.74%        | 77.37 ± 0.58%        | 49.21 ± 2.37%        | 59.97 ± 0.97%        | 21.31 ± 0.77%        |
| AGE+OpenWGL         | 62.27 ± 0.12%        | <u>33.85 ± 0.39%</u> | 82.14 ± 1.04%        | 67.35 ± 3.09%        | 58.89 ± 0.85%        | 22.31 ± 0.85%        |
| AGE+GPN             | 61.34 ± 0.38%        | 31.14 ± 0.82%        | 74.73 ± 1.37%        | 48.62 ± 2.91%        | 56.72 ± 0.73%        | 21.13 ± 0.94%        |
| <b>OWGAL (Ours)</b> | <b>63.32 ± 0.49%</b> | <b>38.23 ± 0.62%</b> | <b>88.41 ± 0.71%</b> | <b>80.12 ± 1.35%</b> | <b>66.46 ± 0.79%</b> | <b>27.15 ± 0.63%</b> |

Better performance on node classification

# Graph Active Learning

## Open world scenario: OWGAL



OWGAL achieves the best performance across different budgets

# Topics for Knowledge from Humans and Domains



Graph learning enhanced by human feedback

- Graph active learning with human knowledge

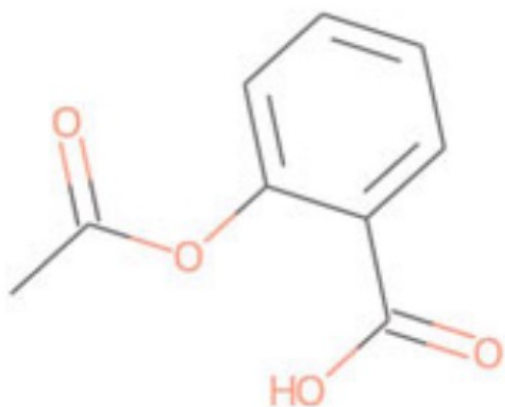
Graph learning enhanced by domain knowledge



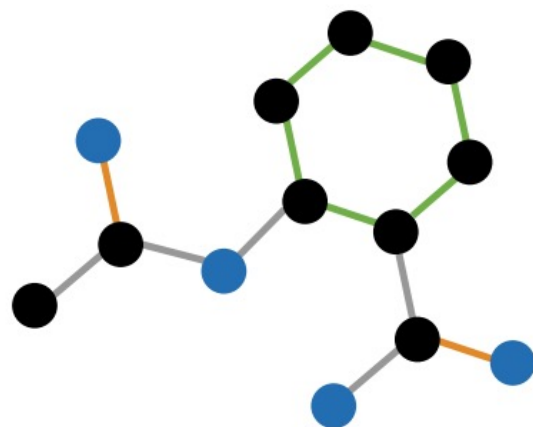
- **Chemistry domain knowledge for molecular property prediction**

# Molecular Graphs

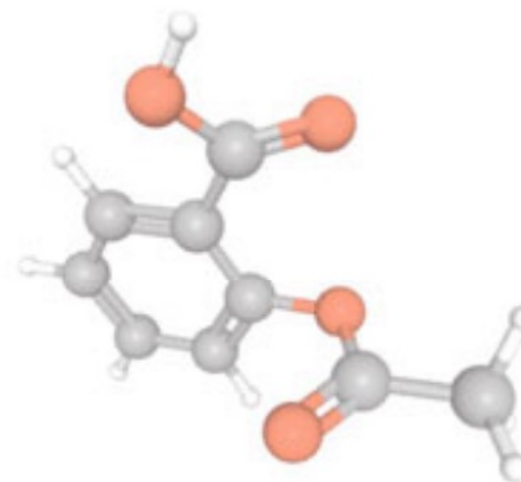
Molecules can be represented as a graph



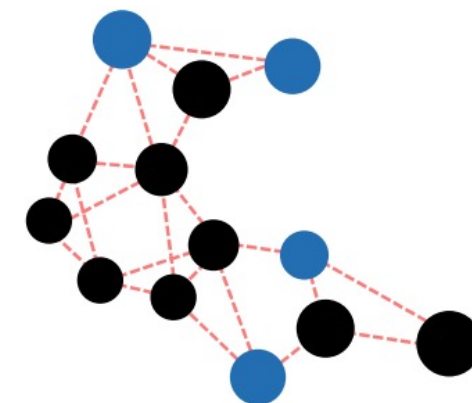
Aspirin in 2D



2D topological graph



Aspirin in 3D

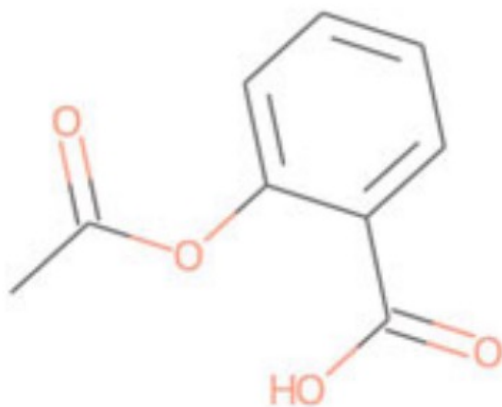


3D geometric graph

# Domain Knowledge for Molecular Property Prediction



## Molecular Property Prediction



model →

Property

Toxicity  
Solubility  
Activity

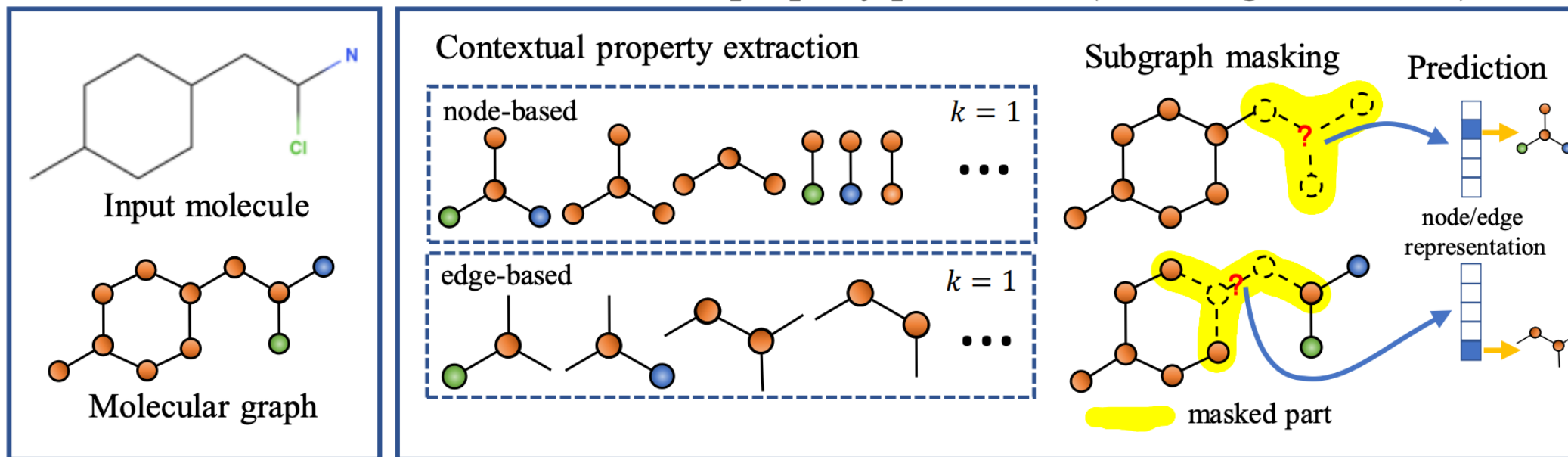
...

# Domain Knowledge for Molecular Property Prediction



Encode molecules in a self-supervised way: GROVER

Contextual property prediction (node/edge level task)



(Task 1) Contextual property prediction: predict the context-aware properties of the target node/edge within some local subgraph

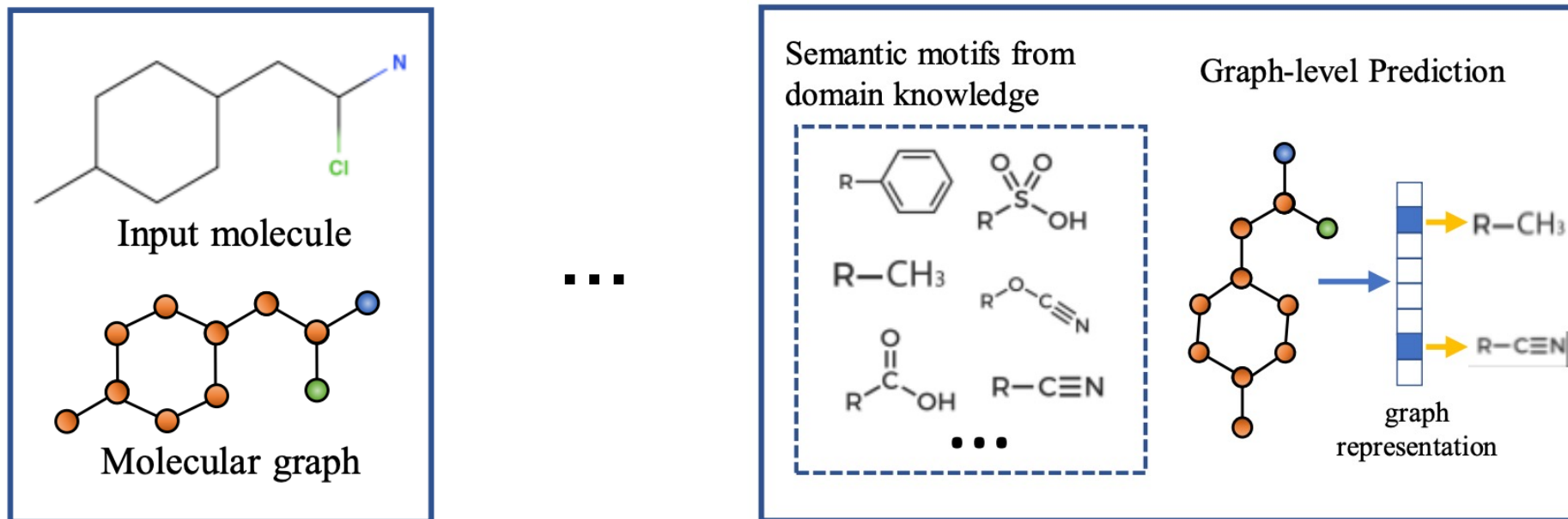


# Domain Knowledge for Molecular Property Prediction

Encode molecules in a self-supervised way: GROVER



## Graph-level motif prediction



(Task 2) Motif prediction: functional groups, as one important class of motifs in molecules, encode the rich domain knowledge of molecules



# Domain Knowledge for Molecular Property Prediction

Encode molecules in a self-supervised way: GROVER



| Classification (Higher is better) |                                 |                                 |                                 |                                 |                                 |                                 |
|-----------------------------------|---------------------------------|---------------------------------|---------------------------------|---------------------------------|---------------------------------|---------------------------------|
| Dataset<br># Molecules            | BBBP<br>2039                    | SIDER<br>1427                   | ClinTox<br>1478                 | BACE<br>1513                    | Tox21<br>7831                   | ToxCast<br>8575                 |
| TF_Robust [40]                    | 0.860 <sub>(0.087)</sub>        | 0.607 <sub>(0.033)</sub>        | 0.765 <sub>(0.085)</sub>        | 0.824 <sub>(0.022)</sub>        | 0.698 <sub>(0.012)</sub>        | 0.585 <sub>(0.031)</sub>        |
| GraphConv [24]                    | 0.877 <sub>(0.036)</sub>        | 0.593 <sub>(0.035)</sub>        | 0.845 <sub>(0.051)</sub>        | 0.854 <sub>(0.011)</sub>        | 0.772 <sub>(0.041)</sub>        | 0.650 <sub>(0.025)</sub>        |
| Weave [23]                        | 0.837 <sub>(0.065)</sub>        | 0.543 <sub>(0.034)</sub>        | 0.823 <sub>(0.023)</sub>        | 0.791 <sub>(0.008)</sub>        | 0.741 <sub>(0.044)</sub>        | 0.678 <sub>(0.024)</sub>        |
| SchNet [45]                       | 0.847 <sub>(0.024)</sub>        | 0.545 <sub>(0.038)</sub>        | 0.717 <sub>(0.042)</sub>        | 0.750 <sub>(0.033)</sub>        | 0.767 <sub>(0.025)</sub>        | 0.679 <sub>(0.021)</sub>        |
| MPNN [13]                         | 0.913 <sub>(0.041)</sub>        | 0.595 <sub>(0.030)</sub>        | 0.879 <sub>(0.054)</sub>        | 0.815 <sub>(0.044)</sub>        | 0.808 <sub>(0.024)</sub>        | 0.691 <sub>(0.013)</sub>        |
| DMPNN [63]                        | 0.919 <sub>(0.030)</sub>        | 0.632 <sub>(0.023)</sub>        | 0.897 <sub>(0.040)</sub>        | 0.852 <sub>(0.053)</sub>        | 0.826 <sub>(0.023)</sub>        | 0.718 <sub>(0.011)</sub>        |
| MGCN [30]                         | 0.850 <sub>(0.064)</sub>        | 0.552 <sub>(0.018)</sub>        | 0.634 <sub>(0.042)</sub>        | 0.734 <sub>(0.030)</sub>        | 0.707 <sub>(0.016)</sub>        | 0.663 <sub>(0.009)</sub>        |
| AttentiveFP [61]                  | 0.908 <sub>(0.050)</sub>        | 0.605 <sub>(0.060)</sub>        | 0.933 <sub>(0.020)</sub>        | 0.863 <sub>(0.015)</sub>        | 0.807 <sub>(0.020)</sub>        | 0.579 <sub>(0.001)</sub>        |
| N-GRAM [29]                       | 0.912 <sub>(0.013)</sub>        | 0.632 <sub>(0.005)</sub>        | 0.855 <sub>(0.037)</sub>        | 0.876 <sub>(0.035)</sub>        | 0.769 <sub>(0.027)</sub>        | - <sup>4</sup>                  |
| HU. et.al[18]                     | 0.915 <sub>(0.040)</sub>        | 0.614 <sub>(0.006)</sub>        | 0.762 <sub>(0.058)</sub>        | 0.851 <sub>(0.027)</sub>        | 0.811 <sub>(0.015)</sub>        | 0.714 <sub>(0.019)</sub>        |
| GROVER <sub>base</sub>            | 0.936 <sub>(0.008)</sub>        | 0.656 <sub>(0.006)</sub>        | 0.925 <sub>(0.013)</sub>        | 0.878 <sub>(0.016)</sub>        | 0.819 <sub>(0.020)</sub>        | 0.723 <sub>(0.010)</sub>        |
| GROVER <sub>large</sub>           | <b>0.940</b> <sub>(0.019)</sub> | <b>0.658</b> <sub>(0.023)</sub> | <b>0.944</b> <sub>(0.021)</sub> | <b>0.894</b> <sub>(0.028)</sub> | <b>0.831</b> <sub>(0.025)</sub> | <b>0.737</b> <sub>(0.010)</sub> |

GROVER achieves the best performance on all datasets compared with existing molecular property predicting methods

# Domain Knowledge for Molecular Property Prediction



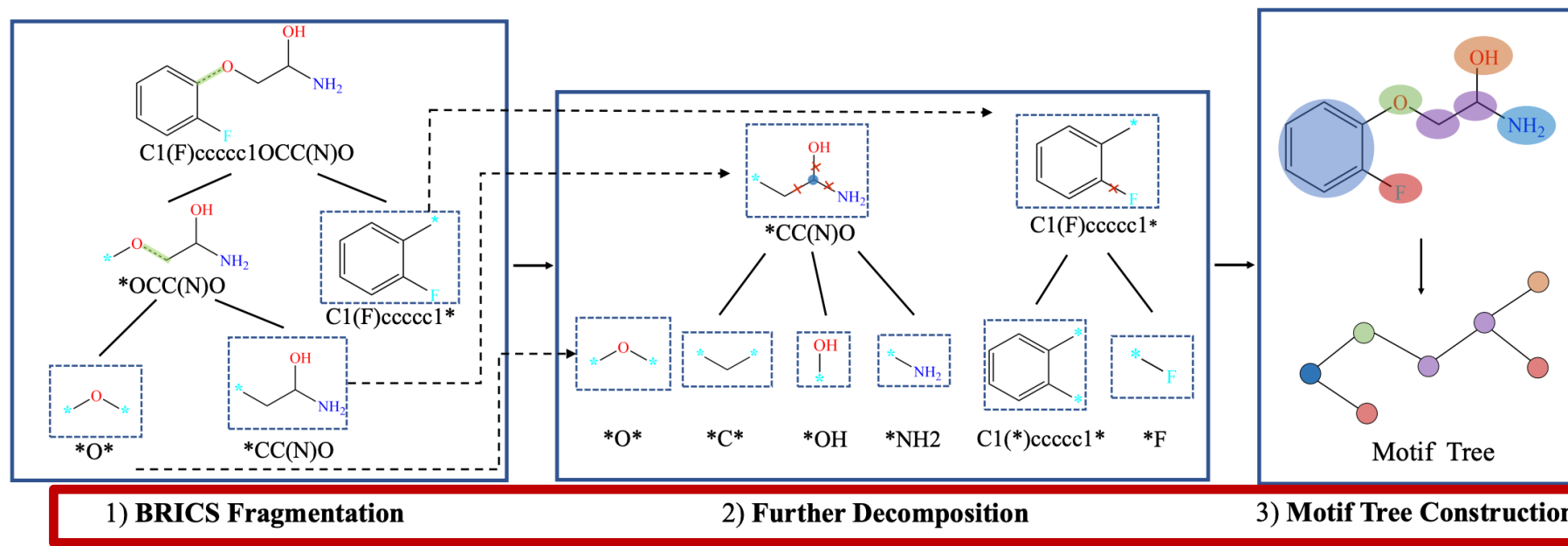
Encode molecules in a self-supervised way: GROVER

|                | GROVER       | No Pretrain | Abs. Imp. |
|----------------|--------------|-------------|-----------|
| BBBP (2039)    | <b>0.940</b> | 0.911       | +0.029    |
| SIDER (1427)   | <b>0.658</b> | 0.624       | +0.034    |
| ClinTox (1478) | <b>0.944</b> | 0.884       | +0.060    |
| BACE (1513)    | <b>0.894</b> | 0.858       | +0.036    |
| Tox21 (7831)   | <b>0.831</b> | 0.803       | +0.028    |
| ToxCast (8575) | <b>0.737</b> | 0.721       | +0.016    |
| Average        | <b>0.834</b> | 0.803       | +0.038    |

Self-supervised pre-training strategy can learn the implicit domain knowledge and enhance the prediction performance

# Domain Knowledge for Molecular Property Prediction

## One step further, generating motifs for molecules: MGSSL



Motif Vocabulary

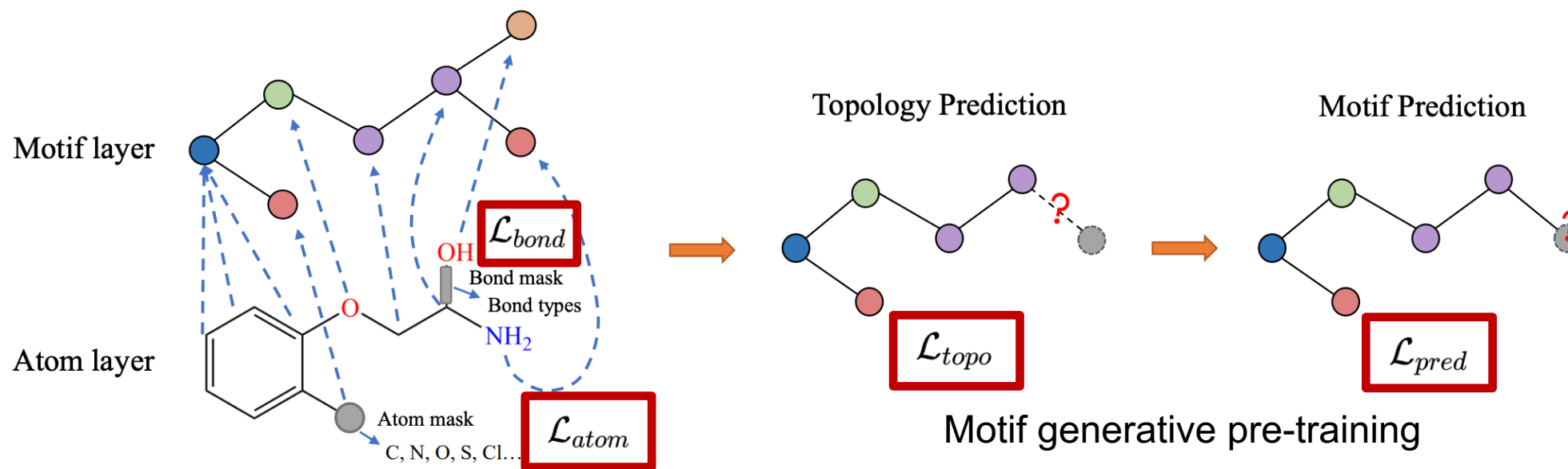


- Three steps: 1) a molecule graph is cleaved based on BRICS. 2) further decomposition to reduce the redundancy of motifs 3) construct motif trees from molecule graphs
- A motif vocabulary can be built via these 3 steps preprocessing

# Domain Knowledge for Molecular Property Prediction



One step further, generating motifs for molecules: MGSSL



Multiple pre-training strategies: 1) reconstruct the masked atom and bond, 2) use the generated motif for motif generative pre-training

# Domain Knowledge for Molecular Property Prediction

One step further, generating motifs for molecules: MGSSL



| SSL methods       | muv             | clintox         | sider           | hiv             | tox21           | bace            | toxcast         | bbbp            | Avg.        |
|-------------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-------------|
| No pretrain       | 71.7±2.3        | 58.2±2.8        | 57.2±0.7        | 75.4±1.5        | 74.3±0.5        | 70.0±2.5        | 63.3±1.5        | 65.5±1.8        | 67.0        |
| Infomax           | 75.1±2.8        | 73.0±3.2        | 58.2±0.5        | 76.5±1.6        | 75.2±0.3        | 75.6±1.0        | 62.8±0.6        | 68.1±1.3        | 70.6        |
| Attribute masking | 74.7±1.9        | 77.5±3.1        | 59.6±0.7        | 77.9±1.2        | <b>77.2±0.4</b> | 78.3±1.1        | 63.3±0.8        | 65.6±0.9        | 71.8        |
| GCC               | 74.1±1.4        | 73.2±2.6        | 58.0±0.9        | 75.5±0.8        | 76.6±0.5        | 75.0±1.5        | 63.5±0.4        | 66.9±0.7        | 70.4        |
| GPT-GNN           | 75.0±2.5        | 74.9±2.7        | 59.3±0.8        | 77.0±1.7        | 76.1±0.4        | 78.5±0.9        | 63.1±0.5        | 67.5±1.3        | 71.4        |
| Grover            | 75.8±1.7        | 76.9±1.9        | 60.7±0.5        | 77.8±1.4        | 76.3±0.6        | 79.5±1.1        | 63.4±0.6        | 68.0±1.5        | 72.3        |
| MGSSL (DFS)       | 78.1±1.8        | 79.7±2.2        | 60.5±0.7        | <b>79.5±1.1</b> | 76.4±0.4        | <b>79.7±0.8</b> | 63.8±0.3        | <b>70.5±1.1</b> | 73.5        |
| MGSSL (BFS)       | <b>78.7±1.5</b> | <b>80.7±2.1</b> | <b>61.8±0.8</b> | 78.8±1.2        | 76.5±0.3        | 79.1±0.9        | <b>64.1±0.7</b> | 69.7±0.9        | <b>73.7</b> |

MGSSL achieves the best performance, demonstrating the effectiveness learning from generated motifs

# Topics for Knowledge from Humans and Domains



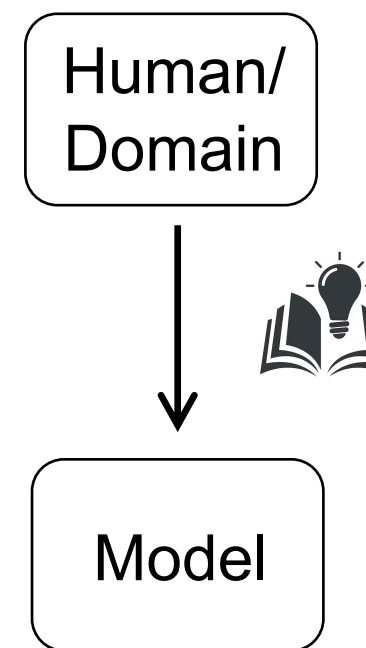
## Chapter summary

### Graph learning enhanced by human feedback

- Graph active learning with human knowledge

### Graph learning enhanced by domain knowledge

- Chemistry domain knowledge for molecular property prediction







AAAI-24 Tutorial  
Feb 2024, Vancouver



## Q & A

Tutorial Website:



[yijuntian.com/tutorial](http://yijuntian.com/tutorial)



# Tutorial Outline



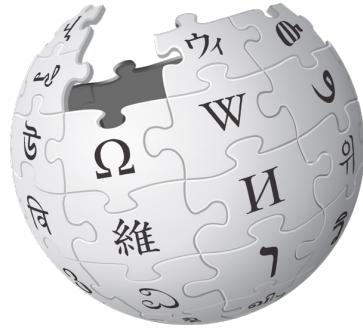
- Preliminaries and Foundations
- Graph Learning Enhanced by Knowledge from Data
- Graph Learning Enhanced by Knowledge from Models
- Graph Learning Enhanced by Knowledge from Humans and Domains
- **Graph Learning Enhanced by Knowledge from External Sources**
- Knowledge-enhanced Graph Learning for Real-world Applications
- Summary and Future Directions



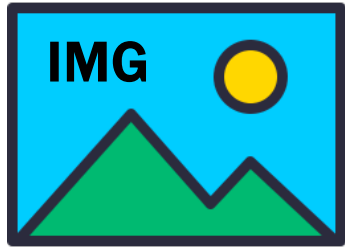
# Knowledge from External Sources



Text



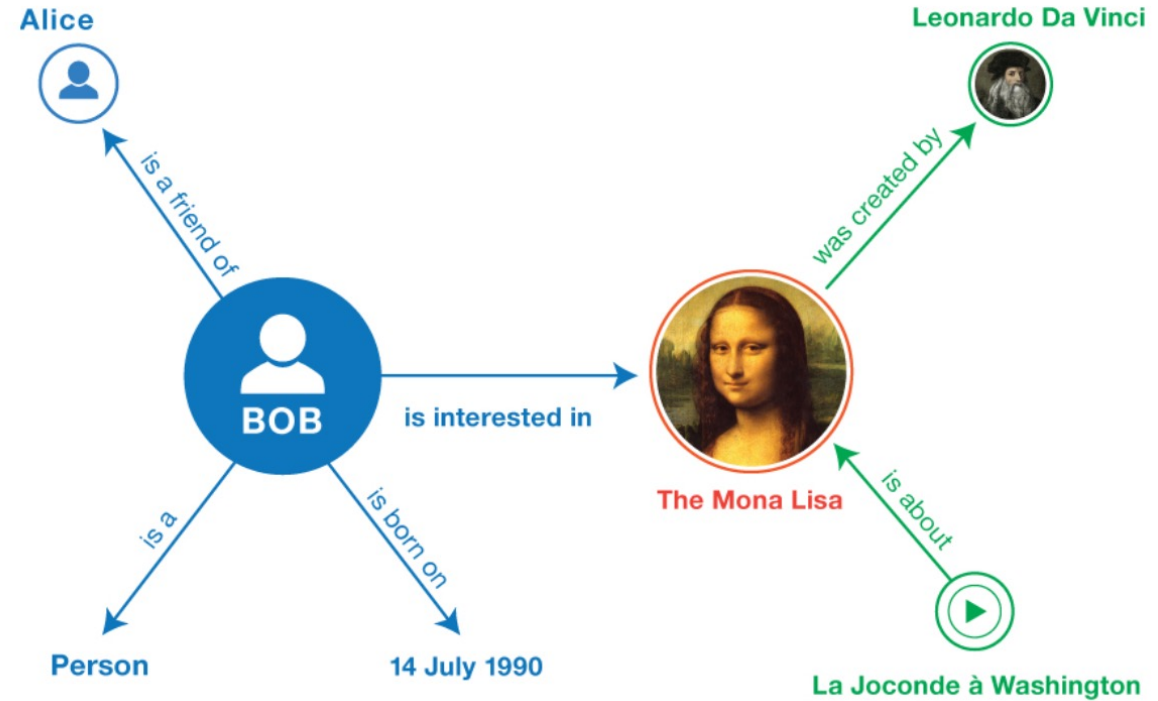
Wikipedia



Image



Video



Knowledge triplets

# Topics for Knowledge from External Sources



## Graph learning on text-rich graphs

- Graph learning on textual-node graphs
- Graph learning on textual-edge graphs

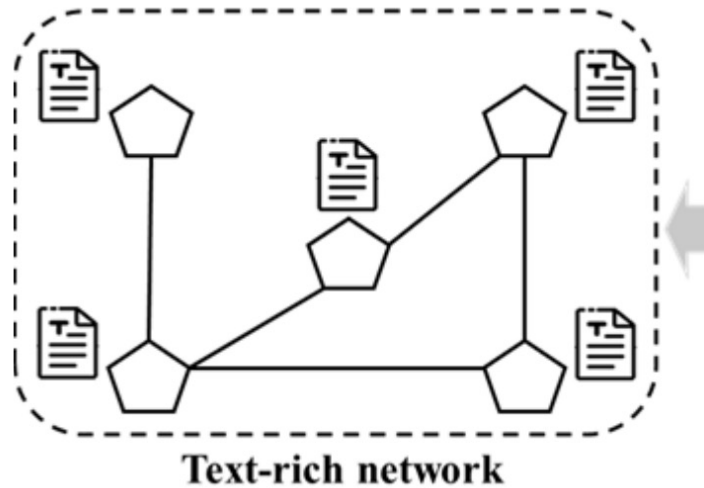
## Graph learning on knowledge graphs

- Knowledge Graph Embedding
- Advancing KG tasks with text data

# Graph Learning on Textual-node Graphs

## What are textual-node graphs?

- A text-rich network contains a node set, an edge set, and a text set
- Each node is associated with some textual information



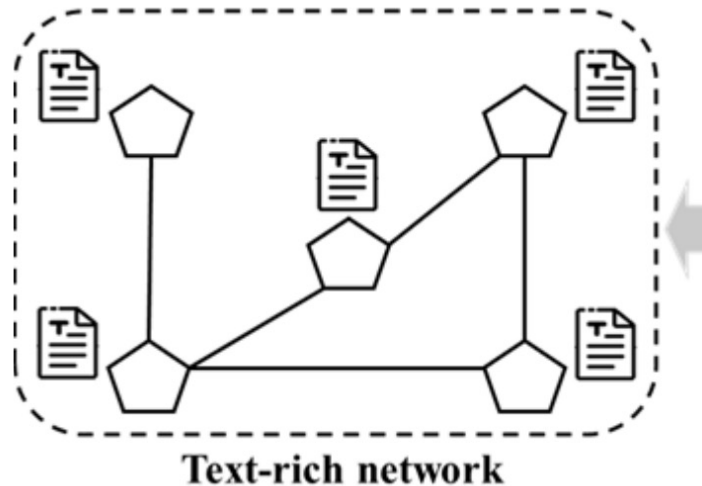
| ID | Document Text                              |
|----|--|
| 1  | Privacy preservation in wireless sensor... |
| 2  | Tabu search for the Steiner problem...     |
| 3  | Feature interaction: a critical review...  |
| 4  | A study of malware in P2P networks...      |
| 5  | Understanding the paradoxical effects...   |

# Graph Learning on Textual-node Graphs



## Why do we need text-rich graphs?

- Text-rich graphs provide additional context and explanation to graphs
- Text allows for modeling the latent correlation between nodes



| ID | Document Text                              |
|----|--|
| 1  | Privacy preservation in wireless sensor... |
| 2  | Tabu search for the Steiner problem...     |
| 3  | Feature interaction: a critical review...  |
| 4  | A study of malware in P2P networks...      |
| 5  | Understanding the paradoxical effects...   |

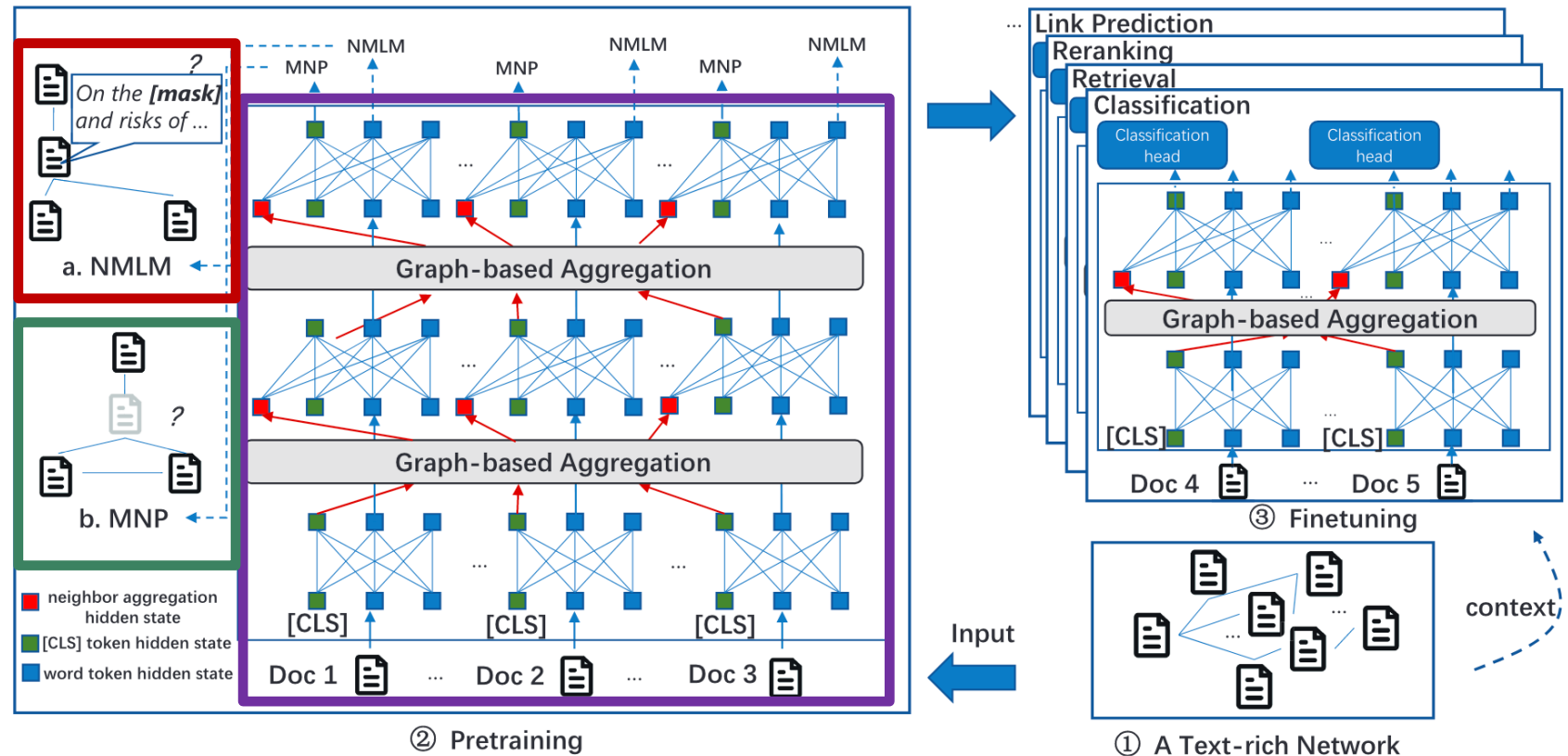
# Graph Learning on Textual-node Graphs

## Pretraining-finetuning on textual-node graphs: Patton



**Network-contextualized Masked Language Modeling:** mask several tokens in the text sequence and utilizes the surrounding unmasked tokens to predict them

**Masked Node Prediction:** Predict the masked nodes based on the adjacent network structure



## Graphformer

# Graph Learning on Textual-node Graphs

## Pretraining-finetuning on text-rich graphs: Patton



Off-the-shelf  
PLMs

Continuous  
pretraining  
method



Consistently  
better

| Method           | Mathematics                  |                              | Geology                      |                              | Economics                    |                              | Clothes                      |                              | Sports                       |                              |
|------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|
|                  | Macro-F1                     | Micro-F1                     | Macro-F1                     | Micro-F1                     | Macro-F1                     | Micro-F1                     | Macro-F1                     | Micro-F1                     | Macro-F1                     | Micro-F1                     |
| BERT             | 18.14 <sub>0.07</sub>        | 22.04 <sub>0.32</sub>        | 21.97 <sub>0.87</sub>        | 29.63 <sub>0.36</sub>        | 14.17 <sub>0.08</sub>        | 19.77 <sub>0.12</sub>        | 45.10 <sub>1.47</sub>        | 68.54 <sub>2.25</sub>        | 31.88 <sub>0.23</sub>        | 34.58 <sub>0.56</sub>        |
| GraphFormers     | 18.69 <sub>0.52</sub>        | 23.24 <sub>0.46</sub>        | 22.64 <sub>0.92</sub>        | 31.02 <sub>1.16</sub>        | 13.68 <sub>1.03</sub>        | 19.00 <sub>1.44</sub>        | 46.27 <sub>1.92</sub>        | 68.97 <sub>2.46</sub>        | 43.77 <sub>0.63</sub>        | 50.47 <sub>0.78</sub>        |
| SciBERT          | 23.50 <sub>0.64</sub>        | 23.10 <sub>2.23</sub>        | 29.49 <sub>1.25</sub>        | 37.82 <sub>1.89</sub>        | 15.91 <sub>0.48</sub>        | 21.32 <sub>0.66</sub>        | -                            | -                            | -                            | -                            |
| SPECTER          | 23.37 <sub>0.07</sub>        | 29.83 <sub>0.96</sub>        | 30.40 <sub>0.48</sub>        | 38.54 <sub>0.77</sub>        | 16.16 <sub>0.17</sub>        | 19.84 <sub>0.47</sub>        | -                            | -                            | -                            | -                            |
| SimCSE (unsup)   | 20.12 <sub>0.08</sub>        | 26.11 <sub>0.39</sub>        | 38.78 <sub>0.19</sub>        | 38.55 <sub>0.17</sub>        | 14.54 <sub>0.26</sub>        | 19.07 <sub>0.43</sub>        | 42.70 <sub>2.32</sub>        | 58.72 <sub>0.34</sub>        | 41.91 <sub>0.85</sub>        | 59.19 <sub>0.55</sub>        |
| SimCSE (sup)     | 20.39 <sub>0.07</sub>        | 25.56 <sub>0.00</sub>        | 25.66 <sub>0.28</sub>        | 33.89 <sub>0.40</sub>        | 15.03 <sub>0.53</sub>        | 18.64 <sub>1.32</sub>        | 52.82 <sub>0.87</sub>        | 75.54 <sub>0.98</sub>        | 46.69 <sub>0.10</sub>        | 59.19 <sub>0.55</sub>        |
| LinkBERT         | 15.78 <sub>0.91</sub>        | 19.75 <sub>1.19</sub>        | 24.08 <sub>0.58</sub>        | 31.32 <sub>0.04</sub>        | 12.71 <sub>0.12</sub>        | 16.39 <sub>0.22</sub>        | 44.94 <sub>2.52</sub>        | 65.33 <sub>4.34</sub>        | 35.60 <sub>0.33</sub>        | 38.30 <sub>0.09</sub>        |
| BERT.MLM         | 23.44 <sub>0.39</sub>        | 31.75 <sub>0.58</sub>        | 36.31 <sub>0.36</sub>        | 48.04 <sub>0.69</sub>        | 16.60 <sub>0.21</sub>        | 22.71 <sub>1.16</sub>        | 46.98 <sub>0.84</sub>        | 68.00 <sub>0.84</sub>        | 62.21 <sub>0.13</sub>        | 75.43 <sub>0.74</sub>        |
| SciBERT.MLM      | 23.34 <sub>0.42</sub>        | 30.11 <sub>0.97</sub>        | 36.94 <sub>0.28</sub>        | 46.54 <sub>0.40</sub>        | 16.28 <sub>0.38</sub>        | 21.41 <sub>0.81</sub>        | -                            | -                            | -                            | -                            |
| SimCSE.in-domain | 25.15 <sub>0.09</sub>        | 29.85 <sub>0.20</sub>        | 38.91 <sub>0.08</sub>        | 48.93 <sub>0.14</sub>        | 18.08 <sub>0.22</sub>        | 23.79 <sub>0.44</sub>        | 57.03 <sub>0.20</sub>        | 80.16 <sub>0.31</sub>        | 65.57 <sub>0.35</sub>        | 75.22 <sub>0.18</sub>        |
| <b>PATTON</b>    | <b>27.58</b> <sub>0.03</sub> | <b>32.82</b> <sub>0.01</sub> | 39.35 <sub>0.06</sub>        | 48.19 <sub>0.15</sub>        | 19.32 <sub>0.05</sub>        | 25.12 <sub>0.05</sub>        | <b>60.14</b> <sub>0.28</sub> | <b>84.88</b> <sub>0.09</sub> | <b>67.57</b> <sub>0.08</sub> | <b>78.60</b> <sub>0.15</sub> |
| SciPATTON        | 27.35 <sub>0.04</sub>        | 31.70 <sub>0.01</sub>        | <b>39.65</b> <sub>0.10</sub> | <b>48.93</b> <sub>0.06</sub> | <b>19.91</b> <sub>0.08</sub> | <b>25.68</b> <sub>0.32</sub> | -                            | -                            | -                            | -                            |
| w/o NMLM         | 25.91 <sub>0.45</sub>        | 27.79 <sub>2.07</sub>        | 38.78 <sub>0.19</sub>        | 48.48 <sub>0.17</sub>        | 18.86 <sub>0.23</sub>        | 24.25 <sub>0.26</sub>        | 56.68 <sub>0.24</sub>        | 80.27 <sub>0.17</sub>        | 65.83 <sub>0.28</sub>        | 76.24 <sub>0.54</sub>        |
| w/o MNP          | 24.79 <sub>0.65</sub>        | 29.44 <sub>1.50</sub>        | 38.00 <sub>0.73</sub>        | 47.82 <sub>1.06</sub>        | 18.69 <sub>0.59</sub>        | 25.63 <sub>1.44</sub>        | 47.35 <sub>1.20</sub>        | 68.50 <sub>2.60</sub>        | 64.23 <sub>1.53</sub>        | 76.03 <sub>1.67</sub>        |

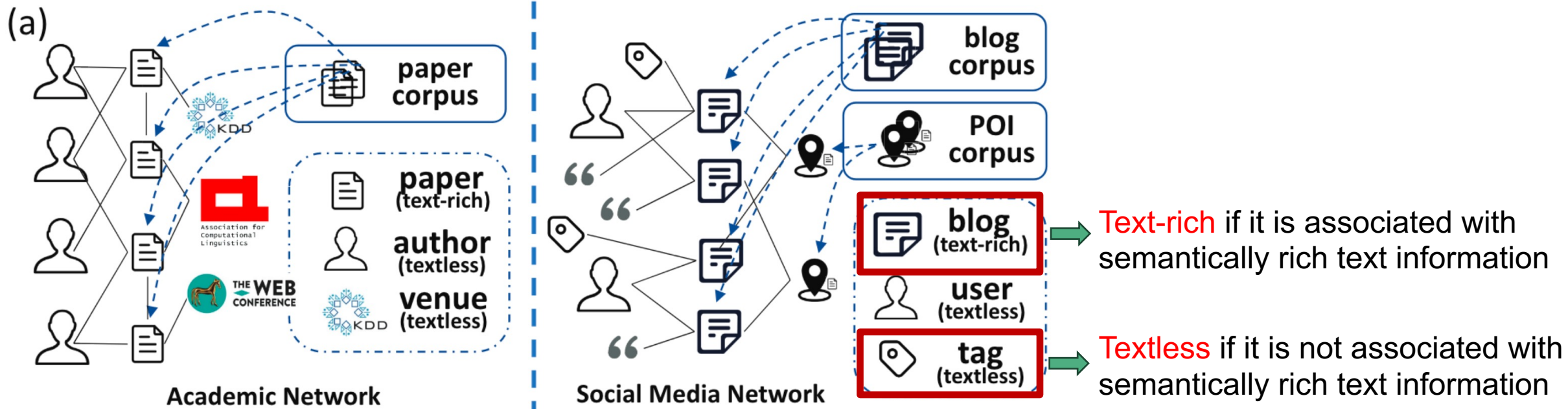


# Graph Learning on Textual-node Graphs



## Heterogeneous text-rich graphs

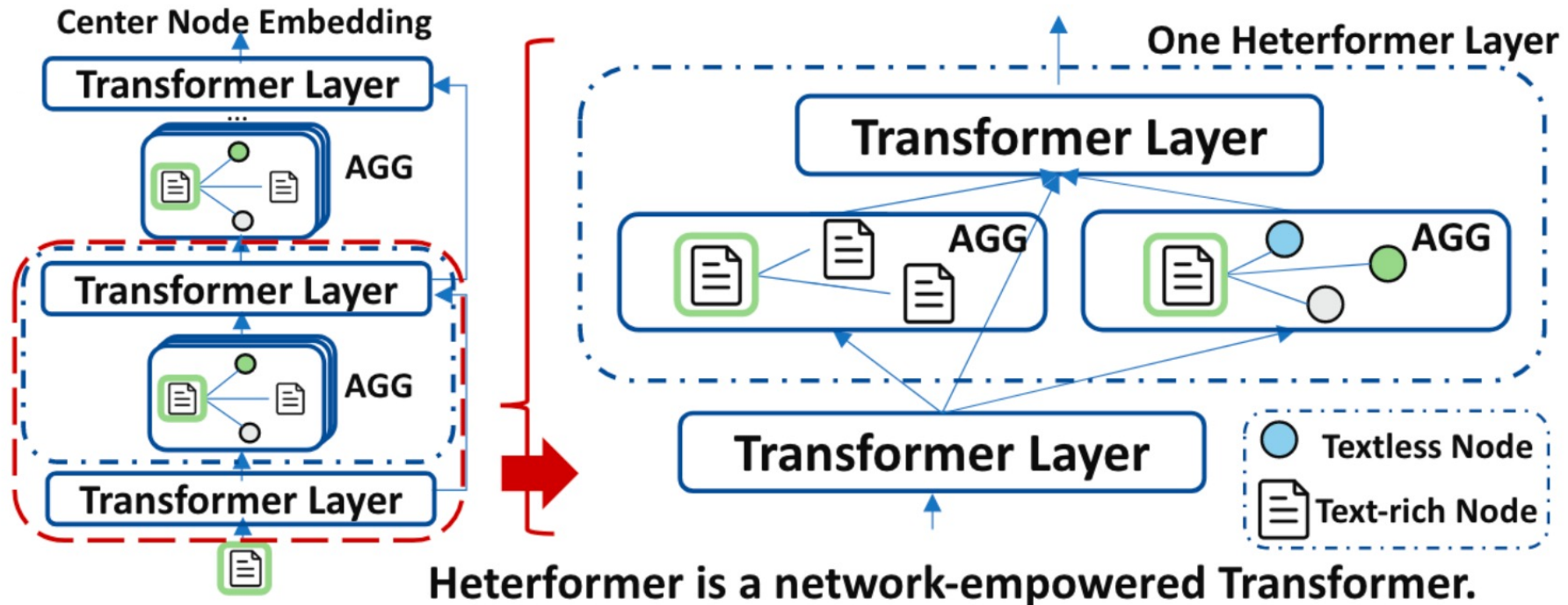
A heterogeneous network includes the sets of nodes, edges, node types, and edge types.



# Graph Learning on Textual-node Graphs

Learn on heterogeneous text-rich networks: Heterformer

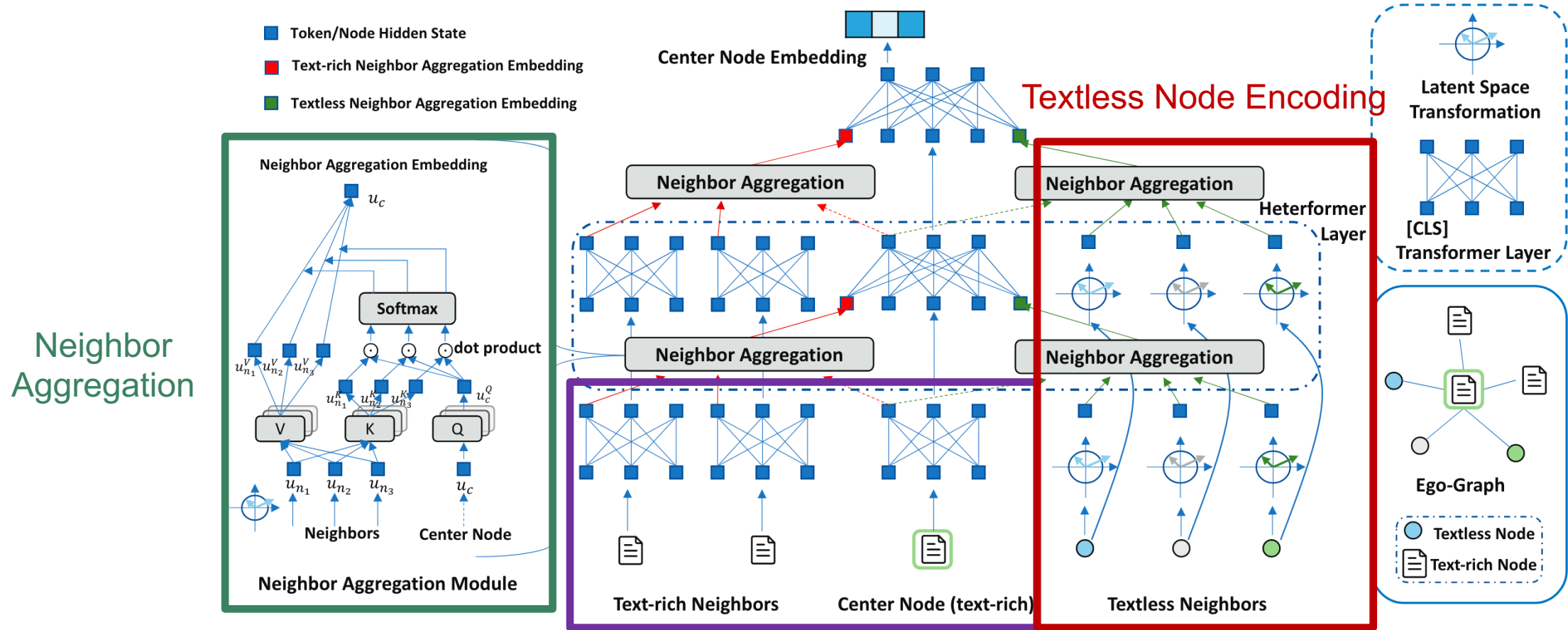
Heterformer is a network-empowered Transformer



# Graph Learning on Textual-node Graphs



## Learn on heterogeneous text-rich networks: Heterformer



Transformer-based Text-Rich Node Encoding

# Graph Learning on Textual-node Graphs

## Learn on heterogeneous text-rich networks: Heterformer



| Method      | DBLP           |                |                | Twitter        |                |                | Goodreads      |                |                |        |
|-------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|--------|
|             | PREC           | MRR            | NDCG           | PREC           | MRR            | NDCG           | PREC           | MRR            | NDCG           |        |
| MeanSAGE    | 0.7019         | 0.7964         | 0.8437         | 0.6489         | 0.7450         | 0.7991         | 0.6302         | 0.7409         | 0.8001         |        |
| BERT        | 0.7569         | 0.8340         | 0.8726         | 0.7179         | 0.7833         | 0.8265         | 0.5571         | 0.6668         | 0.7395         |        |
| Homo GNN    | BERT+MeanSAGE  | 0.8131         | 0.8779         | 0.9070         | 0.7201         | 0.7845         | 0.8275         | 0.7301         | 0.8167         | 0.8594 |
|             | BERT+MAXSAGE   | 0.8193         | 0.8825         | 0.9105         | 0.7198         | 0.7845         | 0.8276         | 0.7280         | 0.8164         | 0.8593 |
|             | BERT+GAT       | 0.8119         | 0.8771         | 0.9063         | 0.7231         | 0.7873         | 0.8300         | 0.7333         | 0.8170         | 0.8593 |
|             | GraphFormers   | 0.8324         | 0.8916         | 0.9175         | 0.7258         | 0.7891         | 0.8312         | 0.7444         | 0.8260         | 0.8665 |
| Hetero GNN  | BERT+RGCN      | 0.7979         | 0.8633         | 0.8945         | 0.7111         | 0.7764         | 0.8209         | 0.7488         | 0.8303         | 0.8699 |
|             | BERT+HAN       | 0.8136         | 0.8782         | 0.9072         | 0.7237         | 0.7880         | 0.8306         | 0.7329         | 0.8174         | 0.8597 |
|             | BERT+HGT       | 0.8170         | 0.8814         | 0.9098         | 0.7153         | 0.7800         | 0.8237         | 0.7224         | 0.8112         | 0.8552 |
|             | BERT+SHGN      | 0.8149         | 0.8785         | 0.9074         | 0.7218         | 0.7866         | 0.8295         | 0.7362         | 0.8195         | 0.8613 |
|             | GraphFormers++ | 0.8233         | 0.8856         | 0.9130         | 0.7159         | 0.7799         | 0.8236         | 0.7536         | 0.8328         | 0.8717 |
| Heterformer | <b>0.8474*</b> | <b>0.9019*</b> | <b>0.9255*</b> | <b>0.7272*</b> | <b>0.7908*</b> | <b>0.8328*</b> | <b>0.7633*</b> | <b>0.8400*</b> | <b>0.8773*</b> |        |

Heterformer captures the heterogeneous structure information and the rich contextualized textual information hidden inside the networks



# Graph Learning on Textual-edge Graphs

## What are textual-edge graphs?

A text-rich network contains a node set, an edge set, and a text set.

Each edge is associated with some textual information

- ★ When a person replies to another on social media, there will be a directed edge between them accompanied by the response texts
- ★ When a user comments on an item, the user's review will be naturally associated with the user-item edge

## Why do we need textual-edge graphs?

- Textual-edge is ubiquitous in the real world
- Textual-edge provides rich information to describe the relationship between nodes

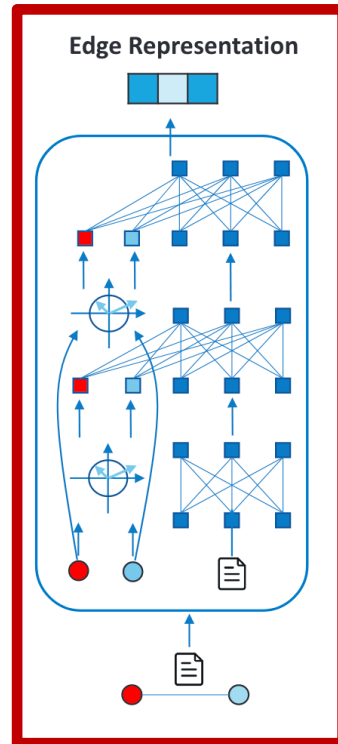


# Graph Learning on Textual-edge Graphs

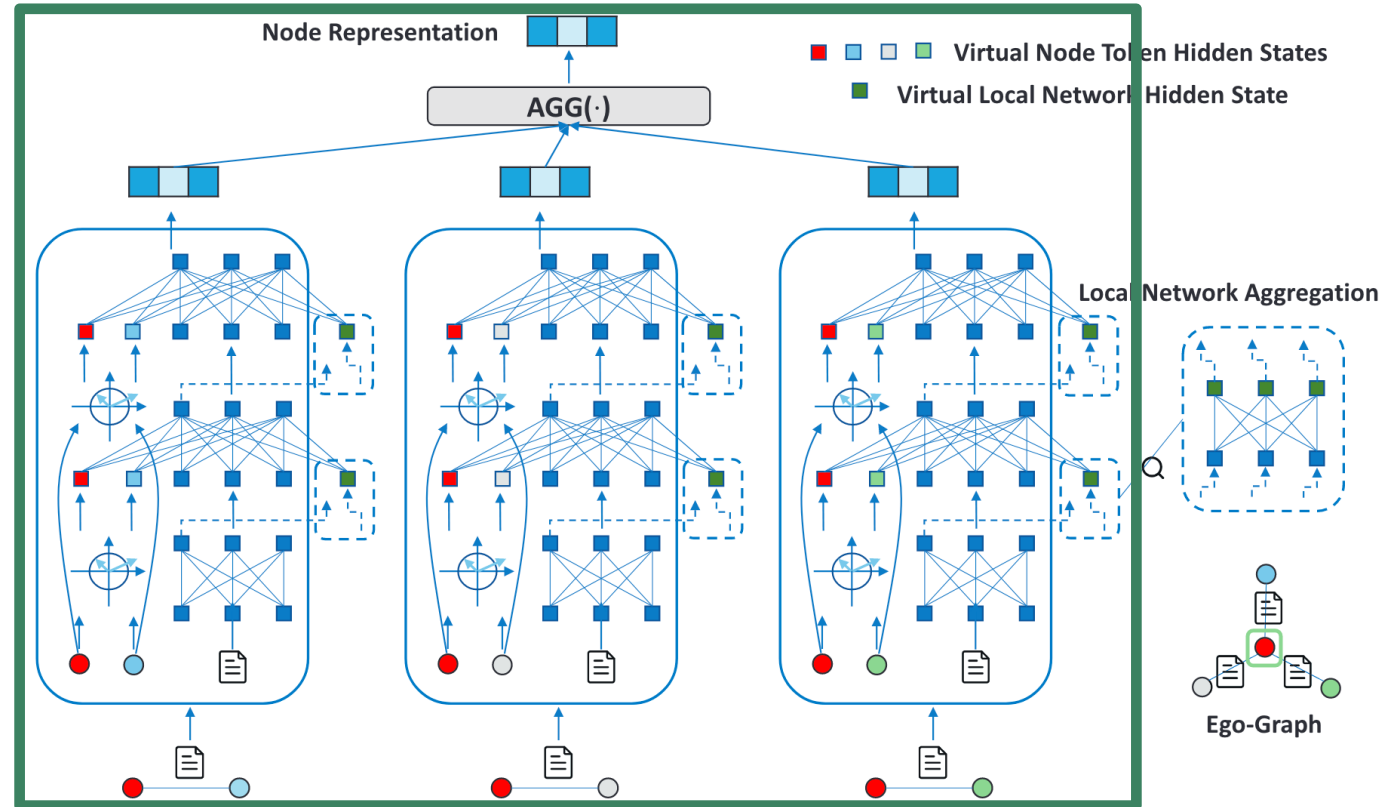
## Learn on textual-edge graphs: Edgeformer



Network-aware  
Edge  
Representation  
Learning



(a) Edgeformer-E



(b) Edgeformer-N

Text-aware Node Representation Learning

# Graph Learning on Textual-edge Graphs



## Learn on textual-edge graphs: Edgeformer

Edge  
Classification

| Model               | Amazon-Movie |              | Amazon-Apps  |              | Goodreads-Crime |              | Goodreads-Children |              |
|---------------------|--------------|--------------|--------------|--------------|-----------------|--------------|--------------------|--------------|
|                     | Macro-F1     | Micro-F1     | Macro-F1     | Micro-F1     | Macro-F1        | Micro-F1     | Macro-F1           | Micro-F1     |
| TF-IDF              | 50.01        | 64.22        | 48.30        | 62.88        | 43.07           | 51.72        | 39.42              | 49.90        |
| TF-IDF+nodes        | 53.59        | 66.34        | 50.56        | 65.08        | 49.35           | 57.50        | 47.32              | 56.78        |
| BERT                | 61.38        | 71.36        | 59.11        | 69.27        | 56.41           | 61.29        | 51.57              | 57.72        |
| BERT+nodes          | 63.00        | 72.45        | 59.72        | 70.82        | 58.64           | 65.02        | 54.42              | 60.46        |
| <b>Edgeformer-E</b> | <b>64.18</b> | <b>73.59</b> | <b>60.67</b> | <b>71.28</b> | <b>61.03</b>    | <b>65.86</b> | <b>57.45</b>       | <b>61.71</b> |

Link  
Prediction

| Model               | Amazon-Movie  |               | Amazon-Apps   |               | Goodreads-Crime |               | Goodreads-Children |               | StackOverflow |               |
|---------------------|---------------|---------------|---------------|---------------|-----------------|---------------|--------------------|---------------|---------------|---------------|
|                     | MRR           | NDCG          | MRR           | NDCG          | MRR             | NDCG          | MRR                | NDCG          | MRR           | NDCG          |
| MF                  | 0.2032        | 0.3546        | 0.1482        | 0.3052        | 0.1923          | 0.3443        | 0.1137             | 0.2716        | 0.1040        | 0.2642        |
| MeanSAGE            | 0.2138        | 0.3657        | 0.1766        | 0.3343        | 0.1832          | 0.3368        | 0.1066             | 0.2647        | 0.1174        | 0.2768        |
| MaxSAGE             | 0.2178        | 0.3694        | 0.1674        | 0.3258        | 0.1846          | 0.3387        | 0.1066             | 0.2647        | 0.1173        | 0.2769        |
| GIN                 | 0.2140        | 0.3648        | 0.1797        | 0.3362        | 0.1846          | 0.3374        | 0.1128             | 0.2700        | 0.1189        | 0.2778        |
| CensNet             | 0.2048        | 0.3568        | 0.1894        | 0.3457        | 0.1880          | 0.3398        | 0.1157             | 0.2726        | 0.1235        | 0.2806        |
| NENN                | 0.2565        | 0.4032        | 0.1996        | 0.3552        | 0.2173          | 0.3670        | 0.1297             | 0.2854        | 0.1257        | 0.2854        |
| BERT                | 0.2391        | 0.3864        | 0.1790        | 0.3350        | 0.1986          | 0.3498        | 0.1274             | 0.2836        | 0.1666        | 0.3252        |
| BERT+MaxSAGE        | 0.2780        | 0.4224        | 0.2055        | 0.3602        | 0.2193          | 0.3694        | 0.1312             | 0.2872        | 0.1681        | 0.3264        |
| BERT+MeanSAGE       | 0.2491        | 0.3972        | 0.1983        | 0.3540        | 0.1952          | 0.3477        | 0.1223             | 0.2791        | 0.1678        | 0.3264        |
| BERT+GIN            | 0.2573        | 0.4037        | 0.2000        | 0.3552        | 0.2007          | 0.3522        | 0.1238             | 0.2801        | 0.1708        | 0.3279        |
| GraphFormers        | 0.2756        | 0.4198        | 0.2066        | 0.3607        | 0.2176          | 0.3684        | 0.1323             | 0.2887        | 0.1693        | 0.3278        |
| BERT+CensNet        | 0.1919        | 0.3462        | 0.1544        | 0.3132        | 0.1437          | 0.3000        | 0.0847             | 0.2436        | 0.1173        | 0.2789        |
| BERT+NENN           | 0.2821        | 0.4256        | 0.2127        | 0.3666        | 0.2262          | 0.3756        | 0.1365             | 0.2925        | 0.1619        | 0.3215        |
| <b>Edgeformer-N</b> | <b>0.2919</b> | <b>0.4344</b> | <b>0.2239</b> | <b>0.3771</b> | <b>0.2395</b>   | <b>0.3875</b> | <b>0.1446</b>      | <b>0.3000</b> | <b>0.1754</b> | <b>0.3339</b> |
| + $\Delta$ %        | 3.5%          | 2.1%          | 5.3%          | 2.9%          | 5.9%            | 3.2%          | 5.9%               | 2.6%          | 2.7%          | 1.8%          |



Significant  
Improvement



# Topics for Knowledge from External Sources



## Graph learning on text-rich graphs

- Graph learning on textual-node graphs
- Graph learning on textual-edge graphs



## Graph learning on knowledge graphs

- Knowledge Graph Embedding
- Advancing KG tasks with text data

# Knowledge Graphs

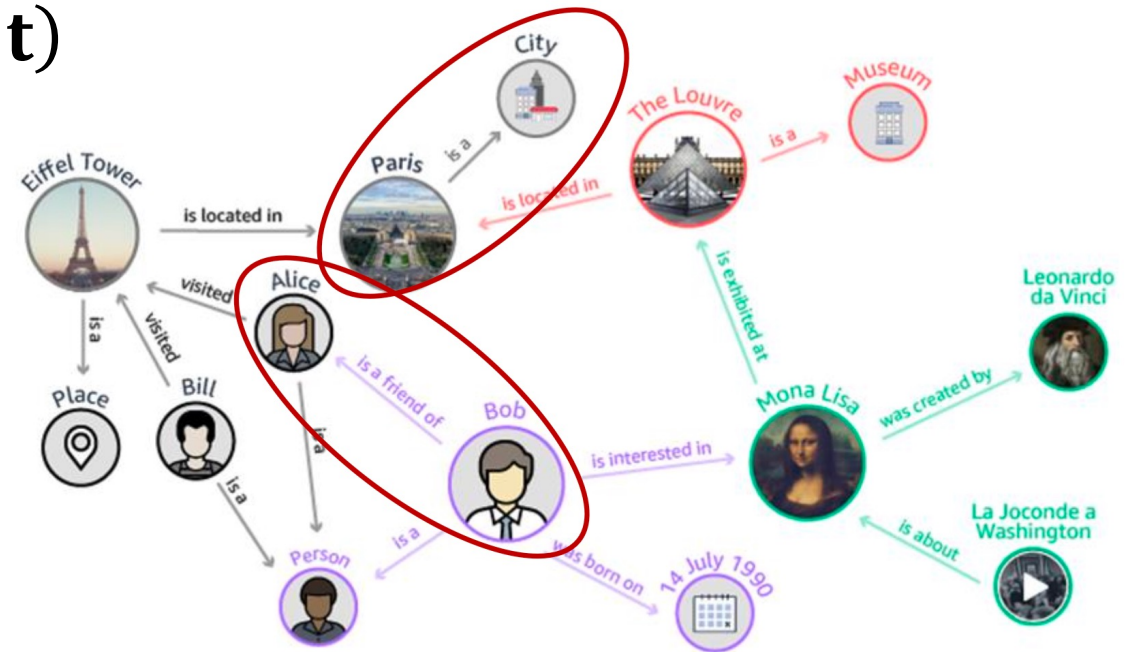
A collection of interlinked entities



- Objects, events or concepts
- Multiple types of entities and relations exist

Facts are represented as triplets (**h**, **r**, **t**)

- ('Paris', 'is\_a', 'city')
- ('Alice', 'is\_friend\_of', 'Bob')



# Knowledge Graph Embedding



## Preliminaries

**Goal:** Encode (1) entities as low-dimensional vectors and (2) relations as parametric algebraic operations in the continuous space

**How-to:** Design a score function  $f_r(\mathbf{h}, \mathbf{t})$  w.r.t. such embedding vectors so that a true triplet receives higher score than a false one

**KGE design rationale:** Capture KG patterns

- Symmetry, antisymmetry, inversion and composition

## Applications of knowledge graph embedding

- Knowledge graph completion
- Question answering
- Recommender system

### Notation & Symbols

- $h$ : head entity
- $r$ : relation
- $t$ : tail entity
- $f_r(\mathbf{h}, \mathbf{t})$ : the score function
- $d_r(\mathbf{h}, \mathbf{t})$ : the distance function
- True/positive triplet:  $(h, r, t)$
- False/negative triplet:  $(h', r, t), (h, r, t'), (h', r, t')$

# Knowledge Graph Embedding

## Preliminaries: Symmetric/Antisymmetric Relations



## Symmetric/Antisymmetric Relations

- **Symmetric:** e.g., Marriage
- **Antisymmetric:** e.g., hasChild

### Formally:

$r$  is **Symmetric:**  $r(x, y) \Rightarrow r(y, x)$  if  $\forall x, y$

$r$  is **Antisymmetric:**  $r(x, y) \Rightarrow \neg r(y, x)$  if  $\forall x, y$

# Knowledge Graph Embedding

## Preliminaries: Inverse Relations



## Inverse Relations

- Hypernym and hyponym:
  - Color is the hypernym ( $r_2$ ) of blue, and blue is the hyponym ( $r_1$ ) of color
- Husband ( $r_2$ ) and wife ( $r_1$ )

## Formally:

$r_1$  is inverse to relation  $r_2$ :  $r_2(x, y) \Rightarrow r_1(y, x)$  if  $\forall x, y$

# Knowledge Graph Embedding

## Preliminaries: Composition Relations



## Composition Relations

- My mother's husband is my father
- $r_1$ : hasMother,  $r_2$ : hasHusband
- $r_3$ : hasFather

## Formally:

$r_3$  is a **composition** of relation  $r_1$  and relation  $r_2$ :

$$r_1(x, y) \wedge r_2(y, z) \Rightarrow r_3(x, z) \text{ if } \forall x, y, z$$



# KG Embedding Method #1: TransE

**Embedding space:** Each entity and relation as a low-dimensional vector in  $R^k$

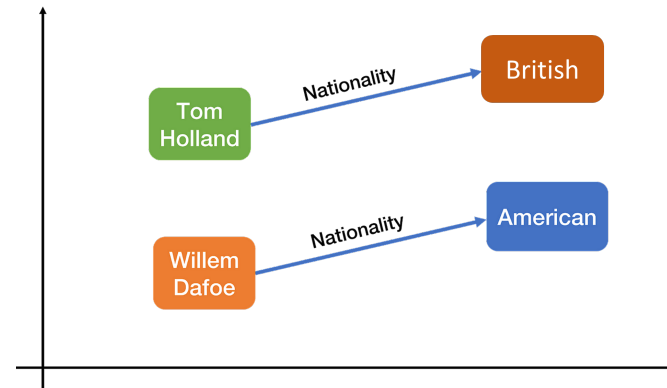
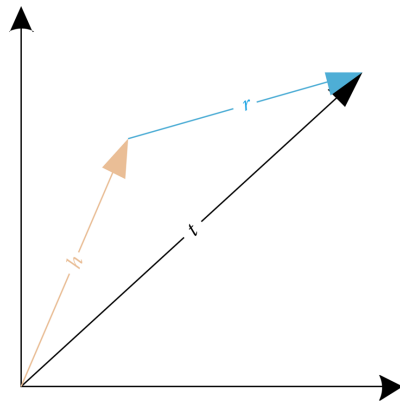
**Key idea:** Relation  $r$  as a translation from the head entity  $h$  to the tail entity  $t$

- An ideal/predicted tail entity:  $\mathbf{t}_{\text{pred}} = \mathbf{h} + \mathbf{r}$

**Score function:**  $f_r(\mathbf{h}, \mathbf{t}) = -d_r(\mathbf{h}, \mathbf{t}) = -\|\mathbf{h} + \mathbf{r} - \mathbf{t}\|$

**Distance function:**  $d_r(\mathbf{h}, \mathbf{t}) = \|\mathbf{h} + \mathbf{r} - \mathbf{t}\|$

- Triplet:  $(h, r, t)$
- Embedding vectors:  $\mathbf{h}, \mathbf{r}, \mathbf{t}$







# KG Embedding Method #1: TransE

## Training Process

For each positive triplet  $(h, r, t) \in S$ ,

- Sample a set of corrupted triplets  $(h, r, t') \in S'_{(h,r,t)}$  or  $(h', r, t) \in S'_{(h,r,t)}$

Learning: Maximize a margin-based ranking criterion

$$\square \text{Minimize } L = \sum_{(h,r,t) \in S} \sum_{(h',r,t') \in S'_{(h,r,t)}} \max(\gamma + d_r(\mathbf{h}, \mathbf{t}) - d_r(\mathbf{h}', \mathbf{t}'), 0)$$

margin

- True triplet:  $(h, r, t)$
- Corrupted triplets:
  - $(h', r, t), (h, r, t'), (h', r, t')$
- $\gamma$ : the margin



# KG Embedding Method #1: TransE

## Key Properties

### Pros

- Can model antisymmetric relations:  $\mathbf{h} + \mathbf{r} = \mathbf{t}$ , but  $\mathbf{t} + \mathbf{r} \neq \mathbf{h}$  if  $\mathbf{r} \neq 0$
- Can model inverse relations:  $\mathbf{h} + \mathbf{r}_1 = \mathbf{t}$ ,  $\mathbf{t} + \mathbf{r}_2 = \mathbf{h}$ ,  $\mathbf{r}_1 = -\mathbf{r}_2$
- Can model composition relations:  $\mathbf{r}_3 = \mathbf{r}_1 + \mathbf{r}_2$

### Cons

- Cannot model symmetric relations:  $\mathbf{h} + \mathbf{r} = \mathbf{t}$ ,  $\mathbf{t} + \mathbf{r} = \mathbf{h}$ , then  $\mathbf{r} = 0$

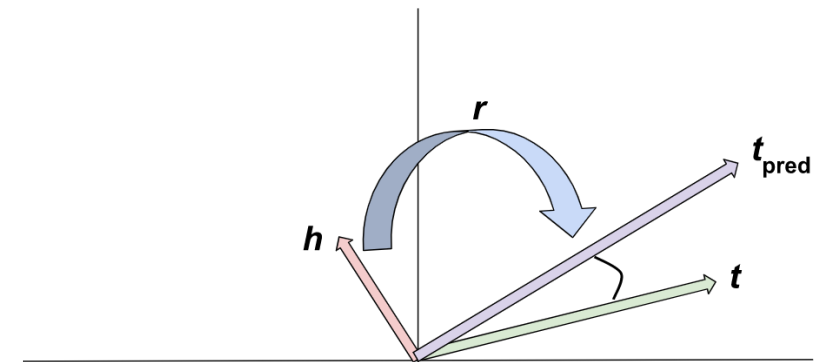
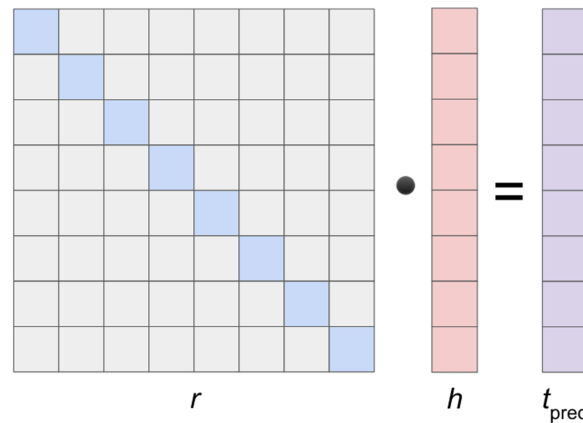
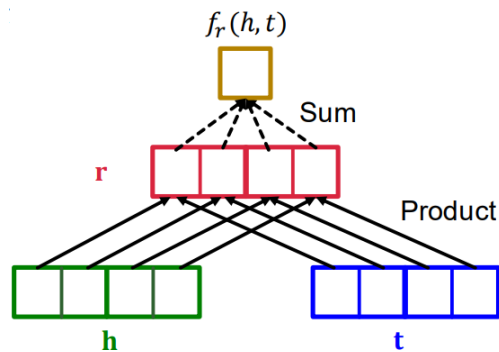
# KG Embedding Method #2: DistMult

**Embedding space:** Each entity and relation as a low-dimensional vector in  $R^k$

**Key idea:** Relation  $r$  defined as the elementwise weights of the head entity

**Score function:**  $f_r(\mathbf{h}, \mathbf{t}) = \langle \mathbf{h}, \mathbf{r}, \mathbf{t} \rangle = \sum_j \mathbf{h}_j \cdot \mathbf{r}_j \cdot \mathbf{t}_j$

- Intuition: Can be viewed as a dot product between  $\mathbf{h} \cdot \mathbf{r}$  and  $\mathbf{t}$



# KG Embedding Method #2: DistMult

## Key Properties

### Pros

- Can model symmetric relations:  $f_r(\mathbf{h}, \mathbf{t}) = \langle \mathbf{h}, \mathbf{r}, \mathbf{t} \rangle = \sum_j \mathbf{h}_j \cdot \mathbf{r}_j \cdot \mathbf{t}_j = \langle \mathbf{t}, \mathbf{r}, \mathbf{h} \rangle = f_r(\mathbf{t}, \mathbf{h})$

### Cons

- Cannot model antisymmetric relations:  $f_r(\mathbf{h}, \mathbf{t})$  and  $f_r(\mathbf{t}, \mathbf{h})$  are always the same
- Cannot model inverse relations: if  $\langle \mathbf{h}, \mathbf{r}_1, \mathbf{t} \rangle = \langle \mathbf{t}, \mathbf{r}_2, \mathbf{h} \rangle$ , it means  $\mathbf{r}_1 = \mathbf{r}_2$
- Cannot model composition relations: it does not model a bijection mapping from  $h$  to  $t$  via relation  $r$

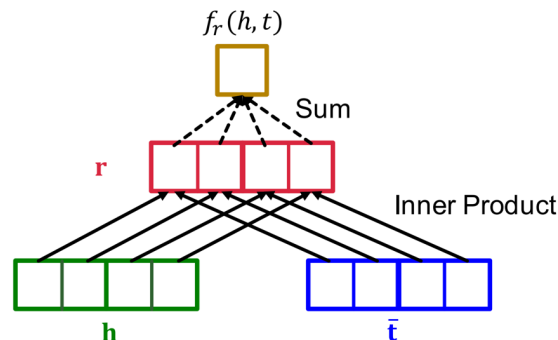
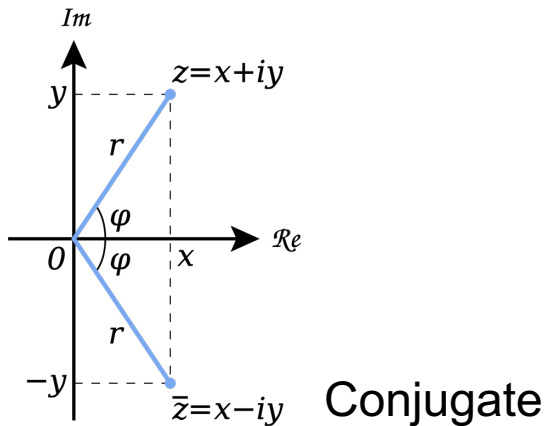


# KG Embedding Method #3: ComplEx

**Embedding space:** Each entity and relation as a low-dimensional vector in  $\mathbb{C}^k$

**Key idea:** Explore the asymmetry of the Hermitian dot product to accommodate antisymmetry

**Score function:**  $f_r(\mathbf{h}, \mathbf{t}) = \text{Re}\left(\sum_j \mathbf{h}_j \cdot \mathbf{r}_j \cdot \bar{\mathbf{t}}_j\right)$



## Notation & Symbols

- $z = x + iy$ : a point in complex space
  - $\text{Re}(z) = x$ : the real part
  - $\text{Im}(z) = y$ : the imaginary part
  - $i = \sqrt{-1}$
- Hermitian dot product  $\langle \mathbf{h}, \bar{\mathbf{t}} \rangle = \sum_j \mathbf{h}_j \cdot \bar{\mathbf{t}}_j$ 
  - $j$ : the  $j$ -th embedding dimension

# KG Embedding Method #3: ComplEx

## Key Properties

### Pros

- Can model antisymmetric relations: when  $\mathbf{Re}(\mathbf{r}) = \mathbf{0}$  (e.g.,  $\mathbf{r}$  is purely imaginary)
- Can model symmetric relations: when  $\mathbf{Im}(\mathbf{r}) = \mathbf{0}$  (e.g.,  $\mathbf{r}$  is purely real)
- Can model inverse: when  $\mathbf{r}_1$  is the conjugate of  $\mathbf{r}_2$

### Cons

- Cannot model composition relations: it does not model a bijection mapping from  $h$  to  $t$  via relation  $r$



# KG Embedding Method #4: RotatE



**Embedding space:** Each entity and relation as a low-dimensional vector in  $\mathbb{C}^k$

**Key idea:** Relations modelled as rotations in complex space

- An ideal tail entity:  $\mathbf{t} = \mathbf{h} \circ \mathbf{r}$ ,

Define each relation  $r$  as an element-wise rotation from the head entity embedding  $\mathbf{h}$  to the tail entity embedding  $\mathbf{t}$ , i.e.,

$$\mathbf{t} = \mathbf{h} \circ \mathbf{r}, \quad \text{where } |\mathbf{r}_j| = 1$$

$\circ$  is the element-wise product. More specifically, we have  $\mathbf{t}_j = \mathbf{h}_j \mathbf{r}_j$ , and

$$\mathbf{r}_j = e^{i\theta_{r,j}}$$

where  $\theta_{r,j}$  is the phase angle of  $r$  in the  $j$ -th dimension

## Notation & Symbols

- $i = \sqrt{-1}$
- $j$ : the  $j$ -th embedding dimension

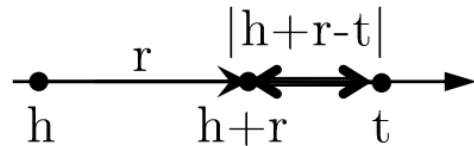


# KG Embedding Method #4: RotatE

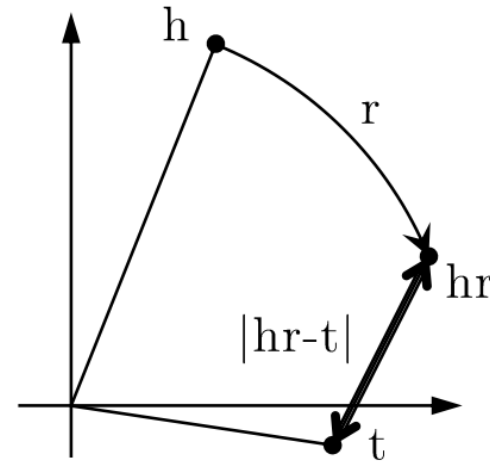
## Geometric Interpretation

**Score function:**  $f_r(\mathbf{h}, \mathbf{t}) = -\|\mathbf{h} \circ \mathbf{r} - \mathbf{t}\|$

Define the distance function of RotatE as  $d_r(\mathbf{h}, \mathbf{t}) = \|\mathbf{h} \circ \mathbf{r} - \mathbf{t}\|$



(a) TransE models  $\mathbf{r}$  as translation in real line.



(b) RotatE models  $\mathbf{r}$  as rotation in complex plane.

# KG Embedding Method #4: RotatE

## Modeling the Relation Patterns

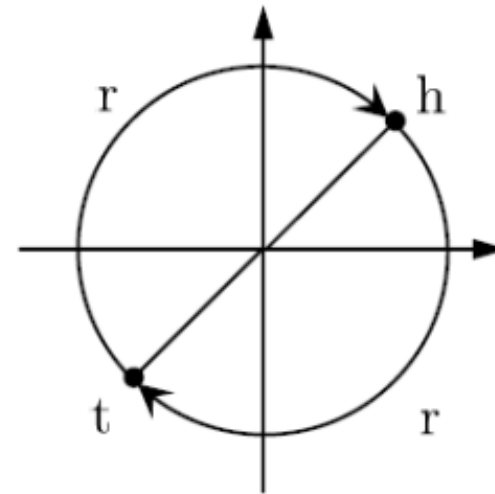
A relation  $r$  is symmetric if and only if  $\mathbf{r}_j = \pm 1$ , i.e.,

$$\theta_{r,j} = 0 \text{ or } \pi$$

$(h, r, t)$  and  $(t, r, h)$

An example on the space of  $\mathcal{C}$

$$\mathbf{r}_j = -1 \text{ or } \theta_{r,j} = \pi$$



# KG Embedding Method #4: RotatE

## Modeling the Relation Patterns



Relation  $r$  is antisymmetric if and only if  $\mathbf{r} \circ \mathbf{r} \neq \mathbf{1}$

Two relations  $r_1$  and  $r_2$  are inverse if and only if  $\mathbf{r}_2 = \bar{\mathbf{r}}_1$ , i.e.,

$$\theta_{2,j} = -\theta_{1,j}$$

A relation  $\mathbf{r}_3 = e^{i\theta_3}$  is a composition of two relations  $\mathbf{r}_1 = e^{i\theta_1}$  and  $\mathbf{r}_2 = e^{i\theta_2}$  if and only if  $\mathbf{r}_3 = \mathbf{r}_1 \circ \mathbf{r}_2$ , i.e.,

$$\theta_3 = \theta_1 + \theta_2$$

# KG Embedding Method #4: RotatE



## Optimization

### Negative sampling loss

- Make the score of true triplet as large as possible, but the score of false one as small as possible

$$L = -\log \sigma(\gamma - d_r(\mathbf{h}, \mathbf{t})) - \sum_{j=1}^k \frac{1}{k} \log \sigma(d_r(\mathbf{h}'_j, \mathbf{t}'_j) - \gamma)$$

Positive triplet      Negative triplet      margin

$\gamma$  is a fixed margin,  $\sigma$  is the sigmoid function, and  $(\mathbf{h}'_j, \mathbf{r}, \mathbf{t}'_j)$  is the  $j$ -th negative triplet

- Controls the minimum separation between the scores of positive and negative triplets.
- Prevent model from over-fitting



# KG Embedding Method #4: RotatE

## Self-adversarial Negative Sampling

Traditionally, the negative samples are drawn in a uniform way

- Inefficient as training goes on since many samples are obviously false
- Does not provide useful information

A self-adversarial negative sampling

- Sample negative triplets according to the current embedding model
- Starts from easier samples to more and more difficult samples
- Curriculum Learning

$$p(h'_j, r, t'_j | \{(h_k, r_k, t_k)\}) = \frac{\exp \alpha f_r(h'_j, t'_j)}{\sum_k \exp \alpha f_r(h'_k, t'_k)}$$

$\alpha$  is the temperature of sampling.  $f_r(h'_j, t'_j)$  measures the salience of the triplet



# KG Embedding Method #4: RotatE

## The Final Objective

Instead of sampling, treating the sampling probabilities as weights

- $p$  is used to weight each negative sample

$$L = -\log \sigma(\gamma - d_r(\mathbf{h}, \mathbf{t})) - \sum_{j=1}^n p(h'_j, r, t'_j) \log \sigma(d_r(\mathbf{h}'_j, \mathbf{t}'_j) - \gamma)$$

Positive triplet      Negative triplet      margin

### Intuitions:

- the smaller the  $d_r(\mathbf{h}'_j, \mathbf{t}'_j)$  is,
- the more likely the triplet is true,
- the harder the negative triplet,
- the higher the weight  $p$

### Notation & Symbols

- $j$ : the  $j$ -th negative sample
- $n$ : the number of negative samples
- $\gamma$ : the margin

# Empirical Comparisons of KG Embedding

## Observation



|          | FB15k-237  |             |             |             |             | WN18RR      |             |             |             |             |
|----------|------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
|          | MR         | MRR         | H@1         | H@3         | H@10        | MR          | MRR         | H@1         | H@3         | H@10        |
| TransE   | 357        | .294        | -           | -           | .465        | 3384        | .226        | -           | -           | .501        |
| DistMult | 254        | .241        | .155        | .263        | .419        | 5110        | .43         | .39         | .44         | .49         |
| Complex  | 339        | .247        | .158        | .275        | .428        | 5261        | .44         | .41         | .46         | .51         |
| ConvE    | 244        | .325        | .237        | .356        | .501        | 4187        | .43         | .40         | .44         | .52         |
| pRotatE  | 178        | .328        | .230        | .365        | .524        | <b>2923</b> | .462        | .417        | .479        | .552        |
| RotatE   | <b>177</b> | <b>.338</b> | <b>.241</b> | <b>.375</b> | <b>.533</b> | 3340        | <b>.476</b> | <b>.428</b> | <b>.492</b> | <b>.571</b> |

RotatE performs the best, on both datasets, on different metrics

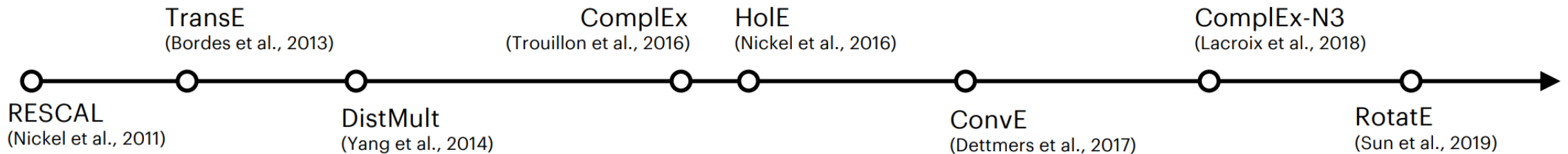


# KG Embedding Method: Summary



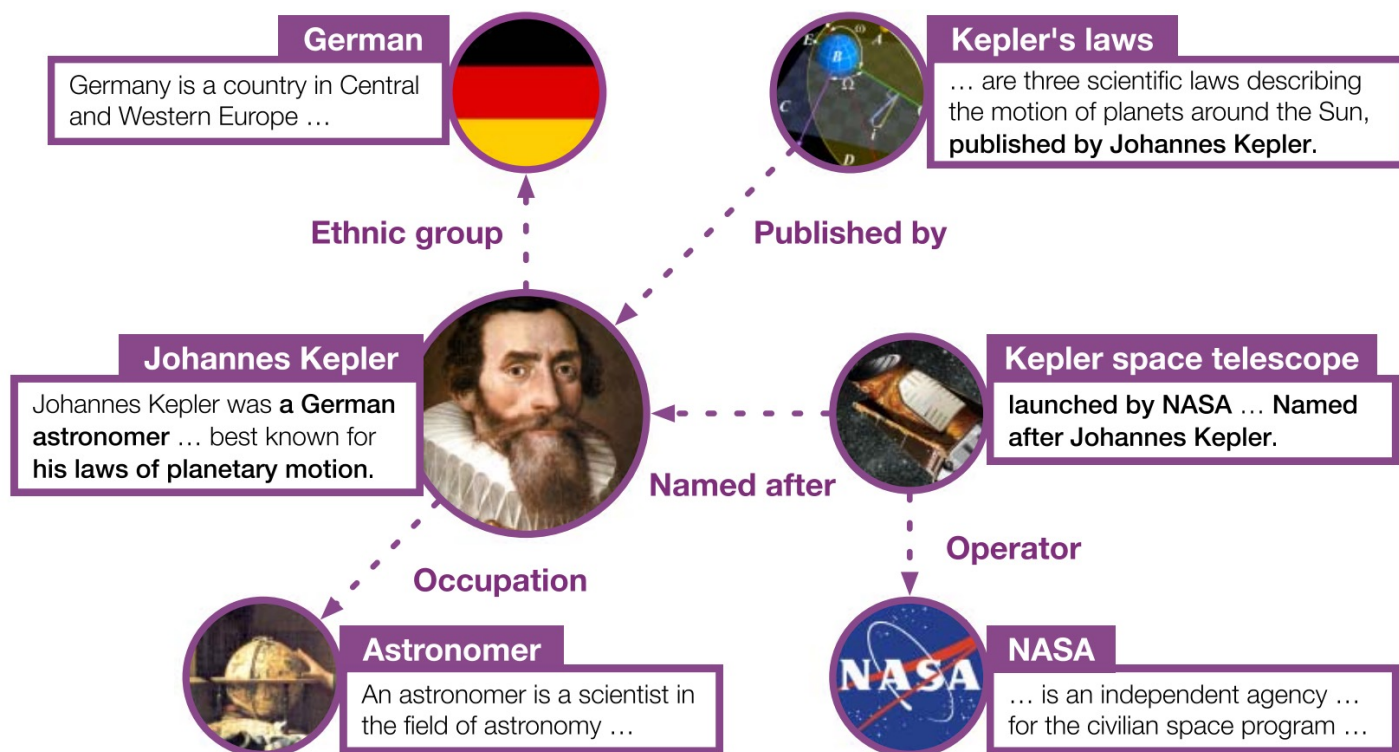
| Model    | Score Function  | Symmetry | Antisymmetry | Inversion | Composition |
|----------|---|----------|--------------|-----------|-------------|
| TransE   | $-  \mathbf{h} + \mathbf{r} - \mathbf{t}  $                     | ✗        | ✓            | ✓         | ✓           |
| DistMult | $\langle \mathbf{h}, \mathbf{r}, \mathbf{t} \rangle$            | ✓        | ✗            | ✗         | ✗           |
| Complex  | $\text{Re}(\langle \mathbf{h}, \mathbf{r}, \mathbf{t} \rangle)$ | ✓        | ✓            | ✓         | ✗           |
| RotatE   | $-  \mathbf{h} \circ \mathbf{r} - \mathbf{t}  $                 | ✓        | ✓            | ✓         | ✓           |

(Some) KGE models in recent literature:



# Advancing KG tasks with text data

## An example of a KG with entity descriptions



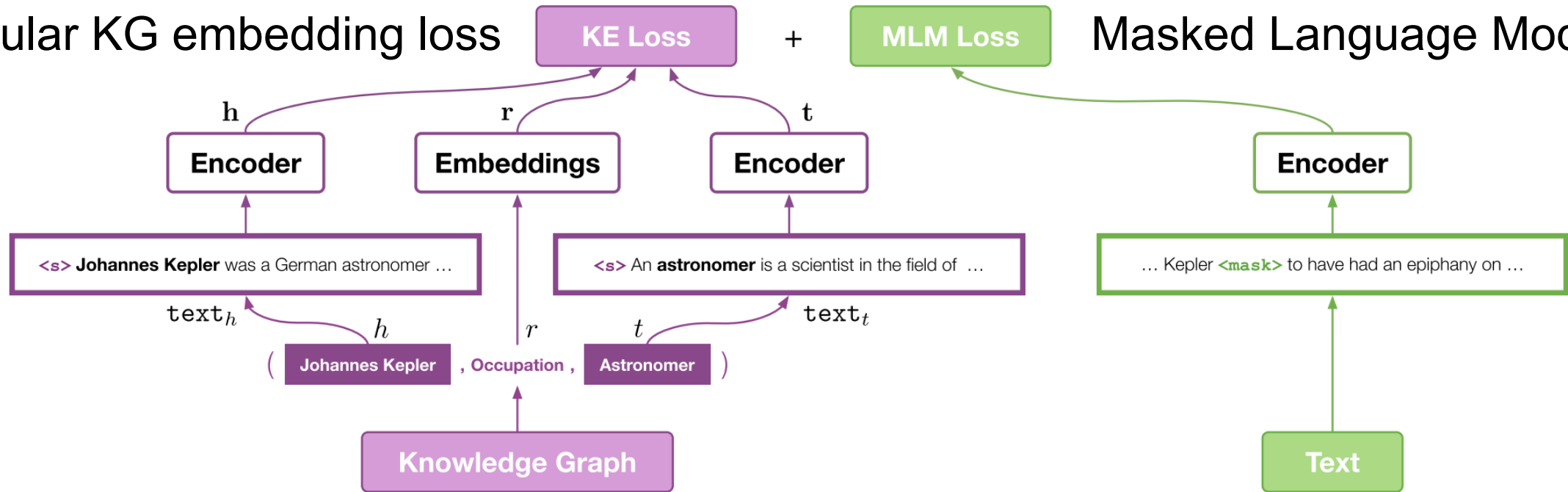
- Descriptions contain abundant information about entities
- Descriptions can help to represent the relational facts between them

# Advancing KG tasks with text data

Train encoder with both KG and NLP losses: KEPLER



Regular KG embedding loss + Masked Language Modeling



KEPLER encodes text descriptions as entity embeddings, and train the shared encoder via MLM loss

# Graph Learning Enhanced by External Knowledge



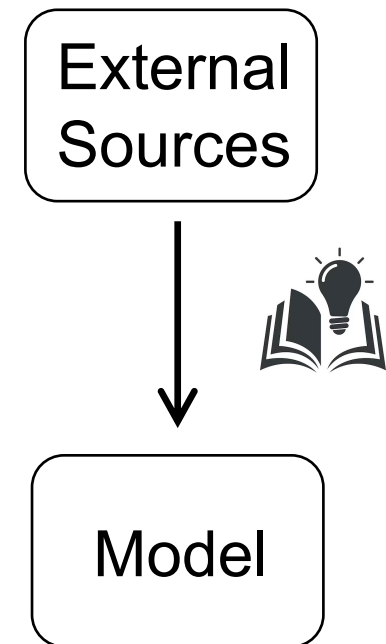
## Chapter summary

### Graph learning on text-rich graphs

- Graph learning on textual-node graphs
- Graph learning on textual-edge graphs

### Graph learning on knowledge graphs

- Knowledge Graph Embedding
- Advancing KG tasks with text data



# Tutorial Outline



- Preliminaries and Foundations
- Graph Learning Enhanced by Knowledge from Data
- Graph Learning Enhanced by Knowledge from Models
- Graph Learning Enhanced by Knowledge from Humans and Domains
- Graph Learning Enhanced by Knowledge from External Sources
- **Knowledge-enhanced Graph Learning for Real-world Applications**
- Summary and Future Directions



# Applications



## 1) Knowledge-enhanced Graph Learning for Recommendation

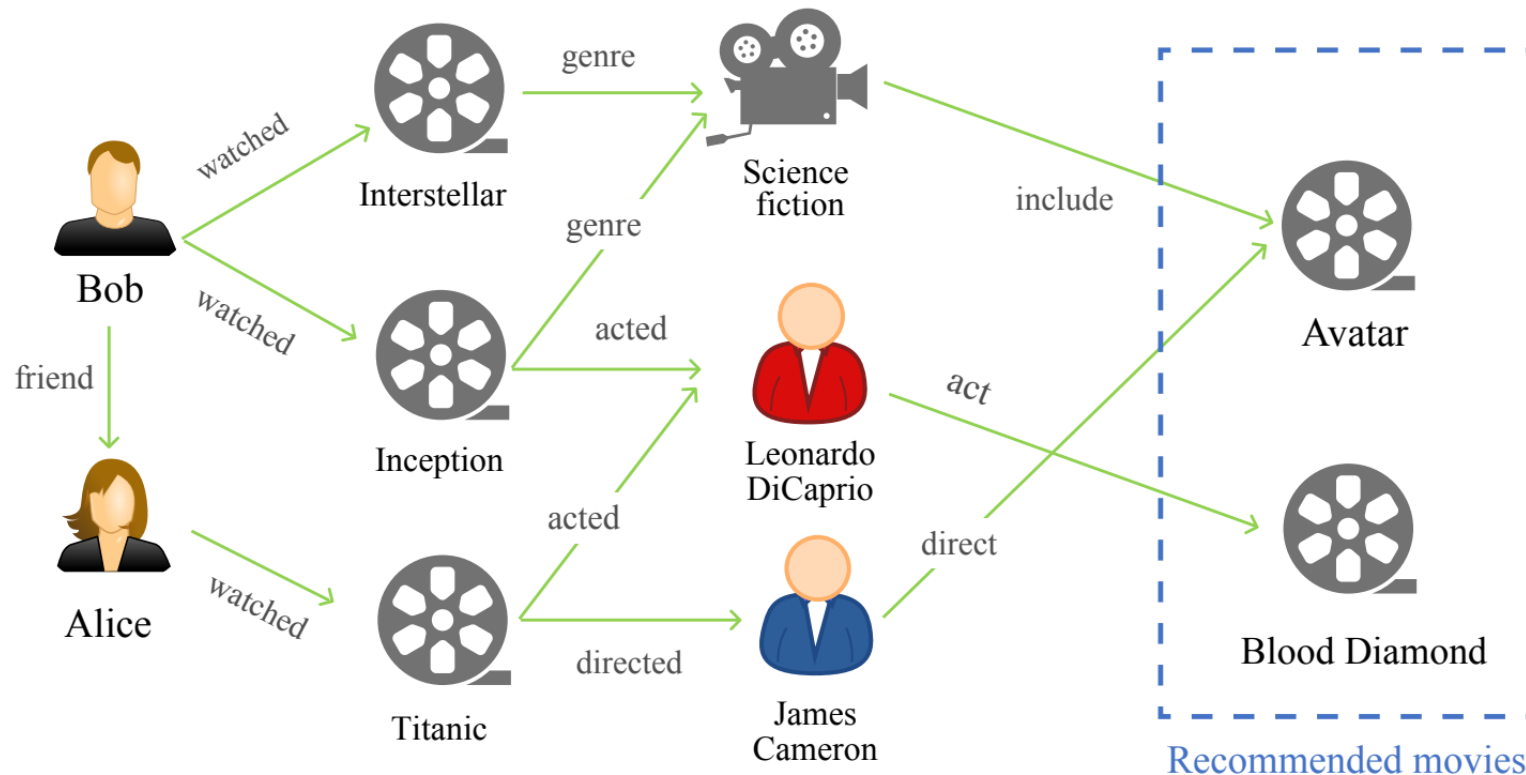
- Path-based recommendation methods
- Propagation-based recommendation methods

## 2) Knowledge-enhanced Graph Learning for Natural Language Processing (NLP)

- Natural Language Understanding
- Commonsense Reasoning
- Advancing LLMs with Knowledge

# Knowledge-enhanced Graph Learning for Recommendation

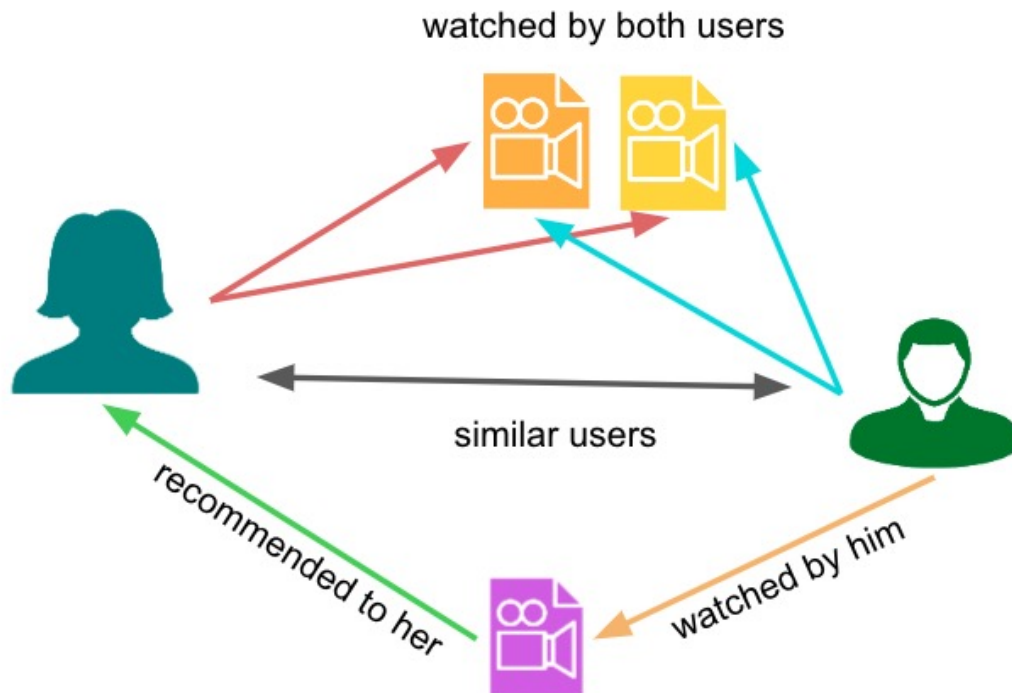
## An example KG for recommendation





# Knowledge-enhanced Graph Learning for Recommendation

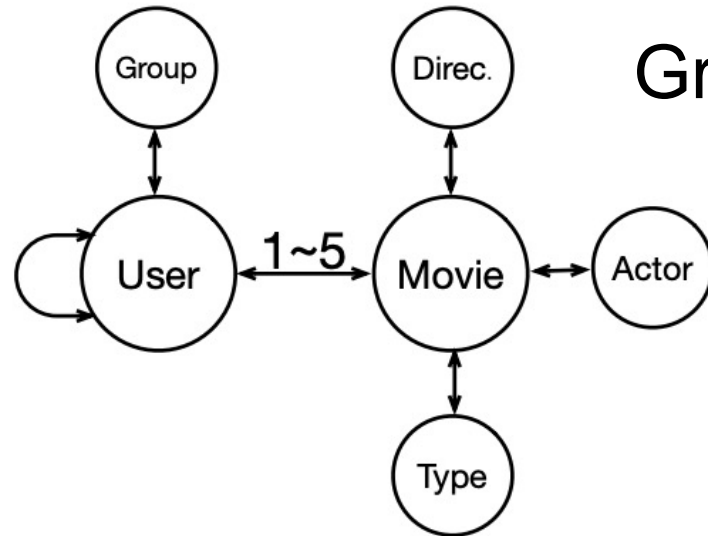
## Preliminary: traditional collaborative filtering



If a person A has the same preference as a person B on a product, A is more likely to have B's preference on a different product than that of a randomly chosen person



## Path-based recommendation methods: SemRec



Graph schema of Douban data

- Metapaths captures different semantic knowledge
- Users appeared in same metapath can have similar ratings

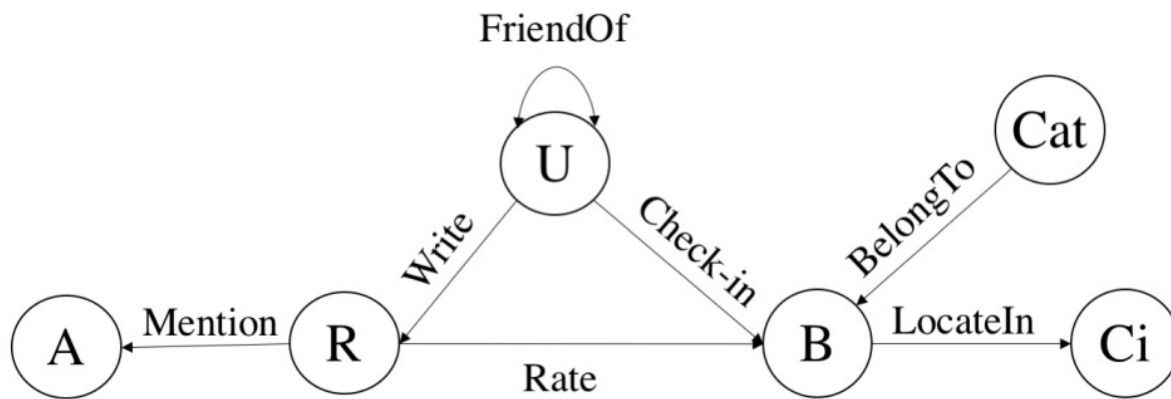
| Meta Path | Semantic Meaning   | Recommendation Model         |
|-----------|--|------------------------------|
| UU        | friends of the target user   | Social recommendation        |
| UGU       | users in the same group of the target user                                   | Member recommendation        |
| UMU       | users who view the same movies with the target user                          | Collaborative recommendation |
| UMTMU     | users who view the movies having the same types with that of the target user | Content recommendation       |

# Knowledge-enhanced Graph Learning for Recommendation



## Problems of SemRec

Metapaths might fail to capture the complex structure and all the semantic knowledge



Graph schema of Yelp data

A: aspect extracted from reviews  
R: reviews  
U: users  
B: business  
Cat: category of item  
Ci: city.

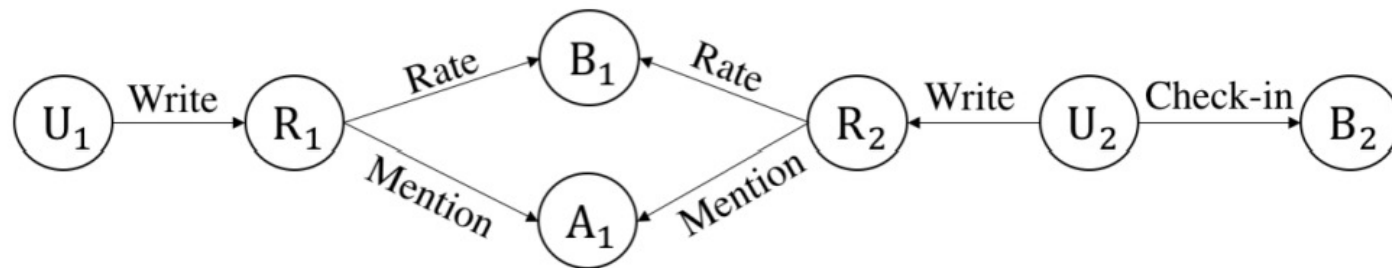
# Knowledge-enhanced Graph Learning for Recommendation



## Address problems of SemRec: FMG

Metapaths might fail to capture the complex structure and all the semantic knowledge

What if  $R_1$  and  $R_2$  give the same rating for a business  $B_1$  and mention the same aspect  $A_1$ ?



Use **metagraph**, a directed acyclic graph that contains more semantic knowledge

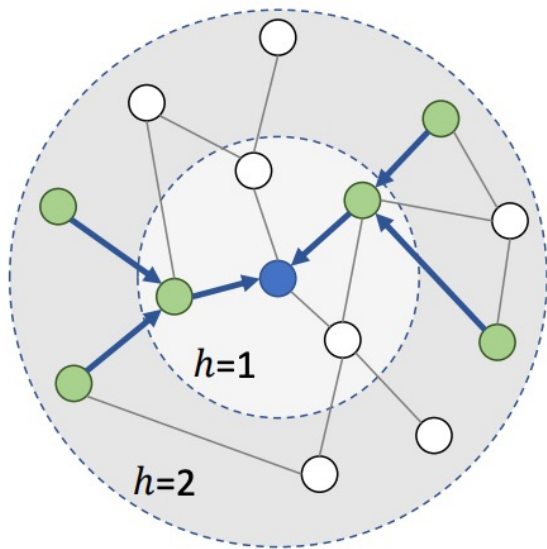


## Path-based recommendation methods: Limitations

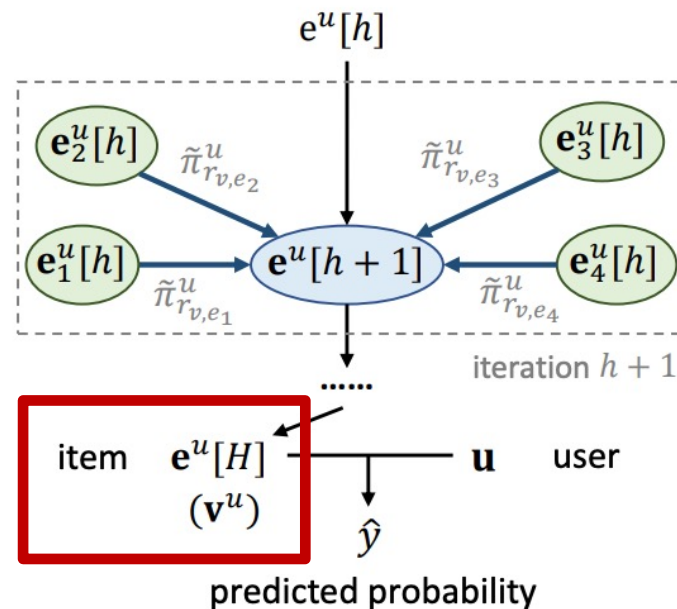
- The choice of metapath and metagraph need to be defined manually
- Metapath and metagraph varying across datasets
- It might require labor-intensive feature engineering to extract relevant paths and structures

# Knowledge-enhanced Graph Learning for Recommendation

## Knowledge graph-aware recommendation methods: KGCN



A two-layer receptive field

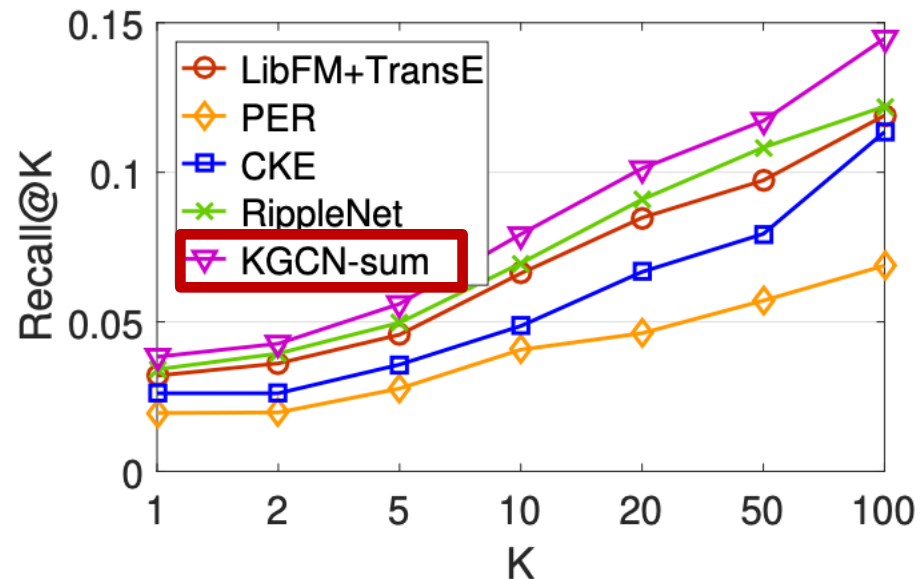


The framework of KGCN

- Learning item representations by aggregation on KG
- Use other items to enhance the representation of target item

# Knowledge-enhanced Graph Learning for Recommendation

## Knowledge graph-aware recommendation methods: KGCN



The results of top-K recommendation on Book-Crossing dataset

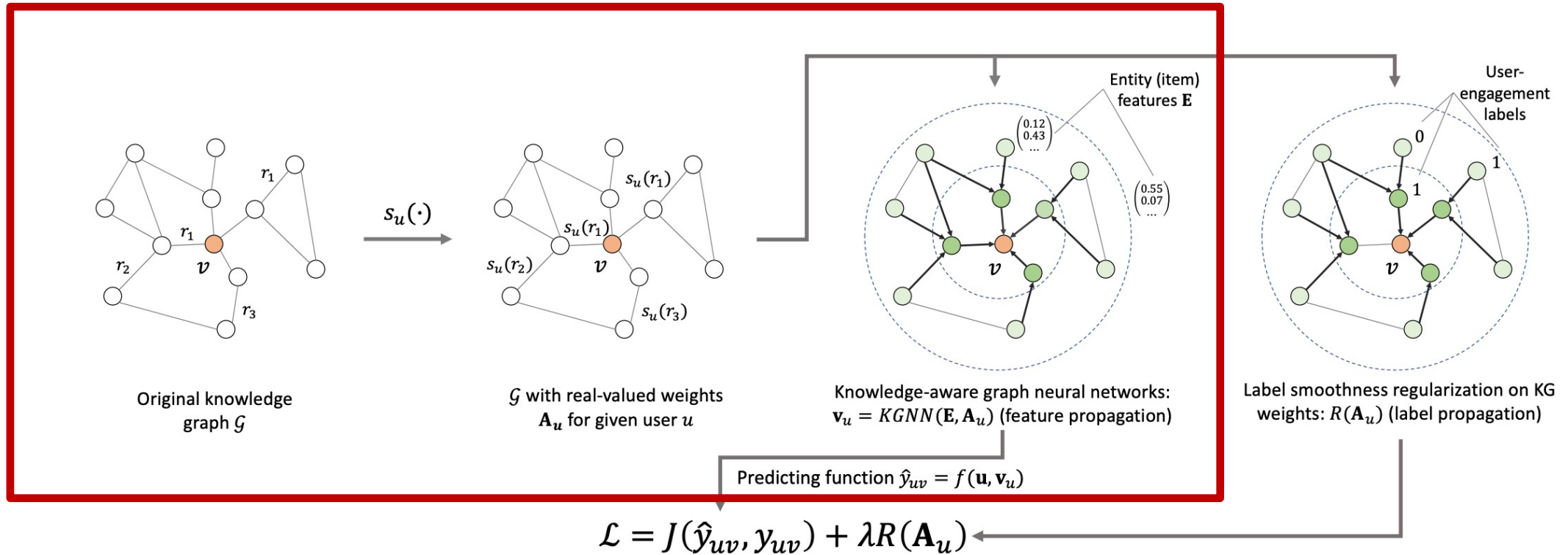
KGCN achieves the best results, demonstrating the benefit of incorporating knowledge from KG



# Knowledge-enhanced Graph Learning for Recommendation



## Knowledge graph-aware recommendation methods: KGNN-LS

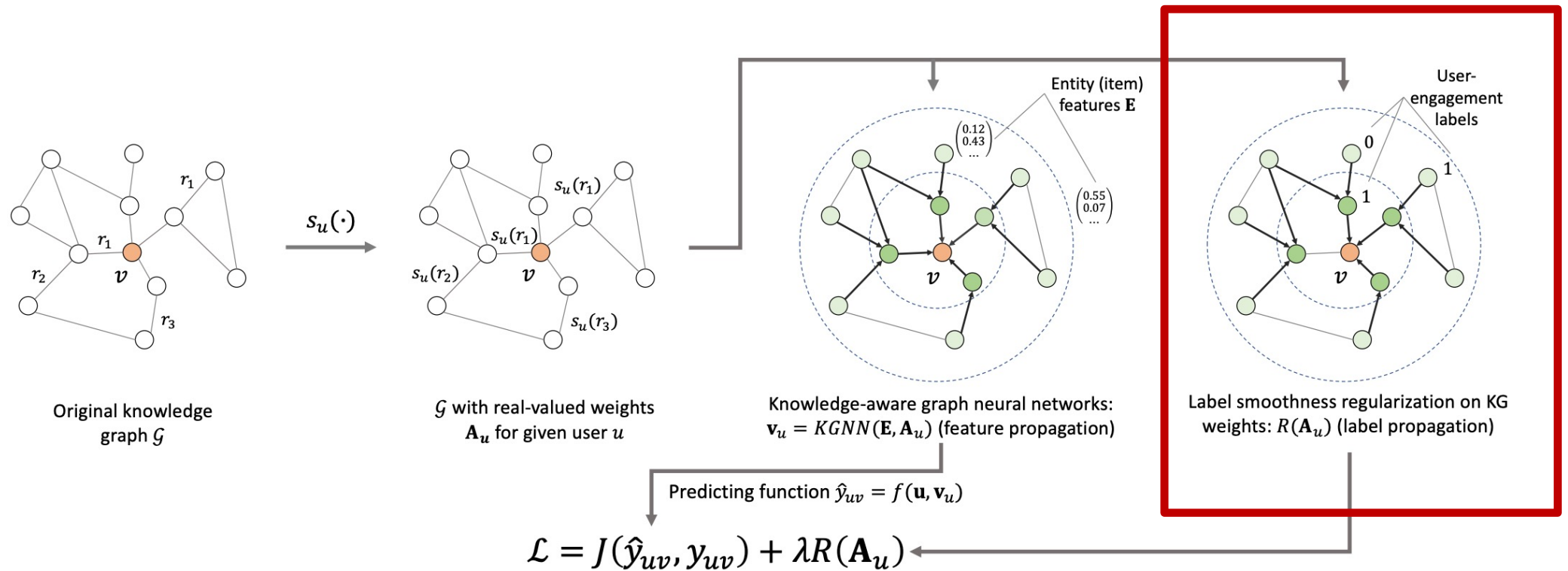


KGNN-LS is Similar to KGCN. First extract a graph of interest to user  $u$ , then use GNNs to learn on this user-specific graph

# Knowledge-enhanced Graph Learning for Recommendation



## Knowledge graph-aware recommendation methods: KGNN-LS

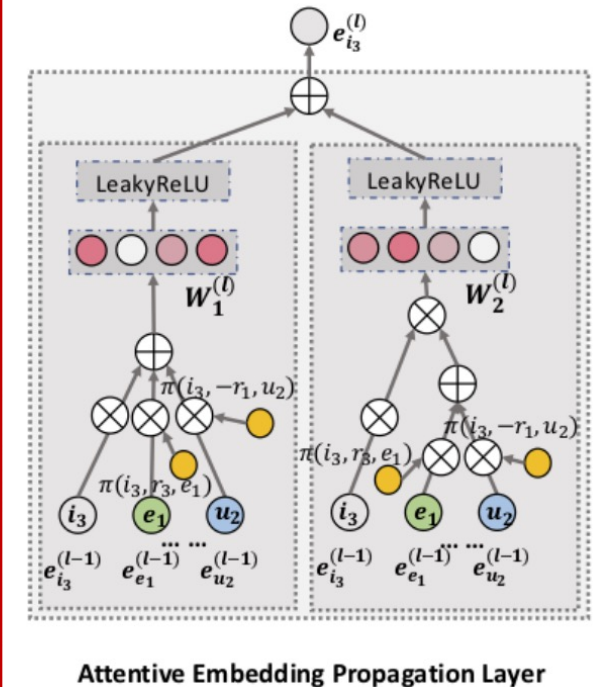
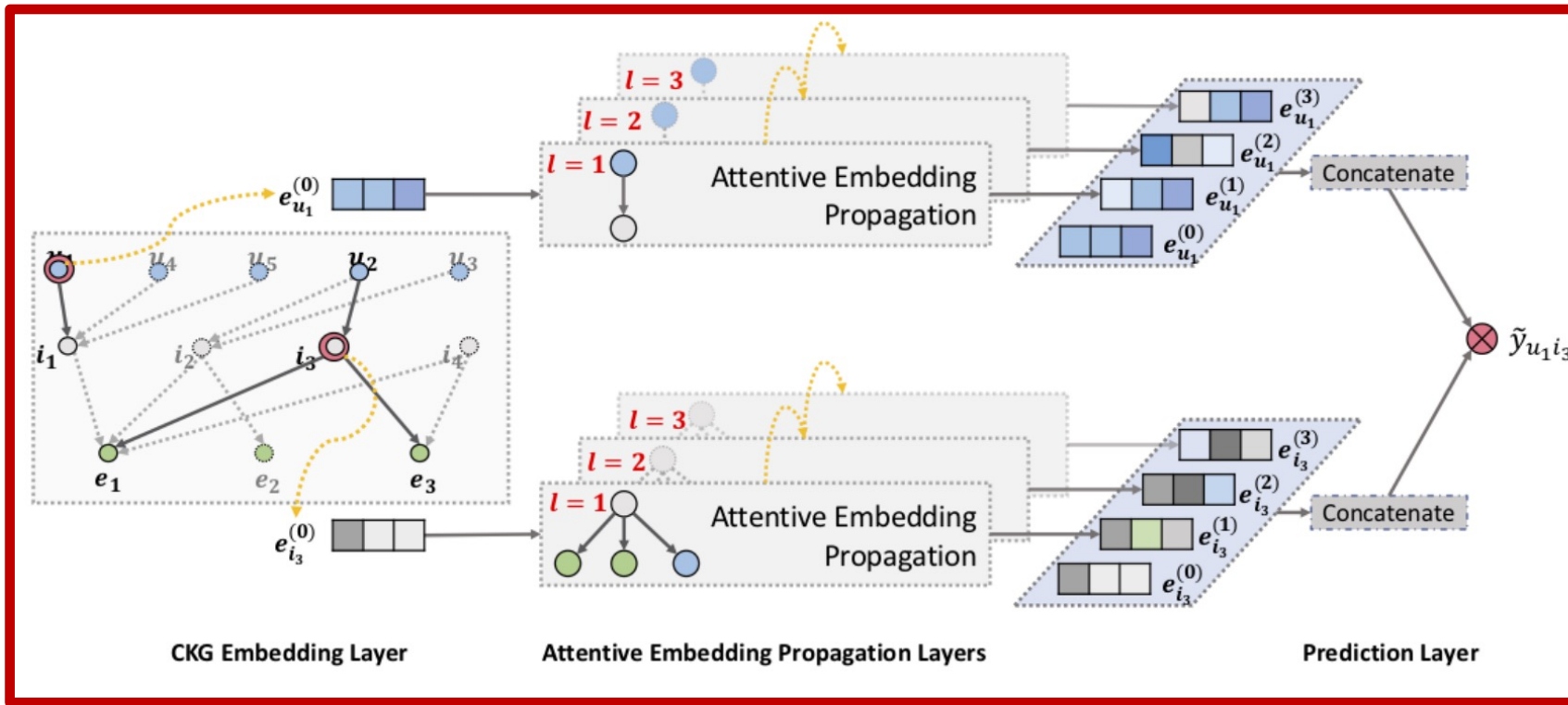


One step further, KGNN-LS utilizes label smoothness: adjacent items in KG are likely to have similar user preferences

# Knowledge-enhanced Graph Learning for Recommendation



## Knowledge graph-aware recommendation methods: KGAT

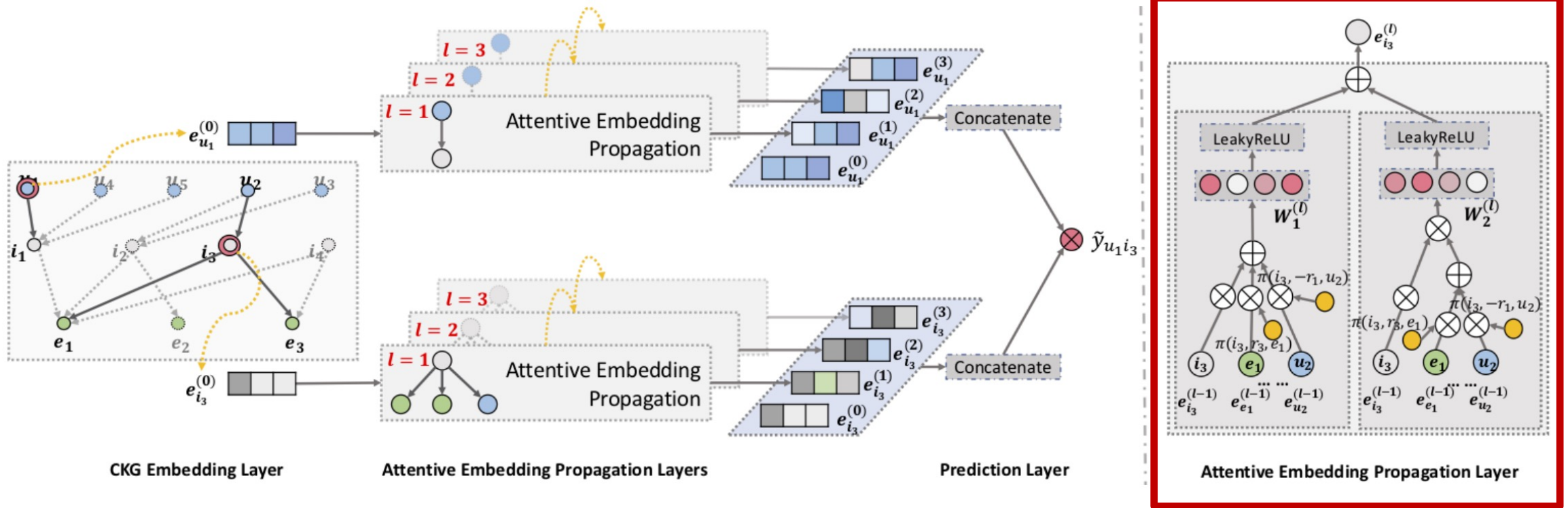


Utilize both item-item connections from KG and user-item direct interaction to learn user and item representations

# Knowledge-enhanced Graph Learning for Recommendation



## Knowledge graph-aware recommendation methods: KGAT



Aggregate neighbor information with attention mechanism: reveal the importance of high-order connections



# Knowledge-enhanced Graph Learning for Recommendation

## Knowledge graph-aware recommendation methods: KGAT



Recommendation Performance Comparison

|           | Amazon-Book    |                | Last-FM        |                | Yelp2018       |                |
|-----------|----------------|----------------|----------------|----------------|----------------|----------------|
|           | recall         | ndcg           | recall         | ndcg           | recall         | ndcg           |
| FM        | 0.1345         | 0.0886         | 0.0778         | 0.1181         | 0.0627         | 0.0768         |
| NFM       | <b>0.1366</b>  | 0.0913         | <b>0.0829</b>  | 0.1214         | 0.0660         | 0.0810         |
| CKE       | 0.1343         | 0.0885         | 0.0736         | 0.1184         | 0.0657         | 0.0805         |
| CFKG      | 0.1142         | 0.0770         | 0.0723         | 0.1143         | 0.0522         | 0.0644         |
| MCTRec    | 0.1113         | 0.0783         | -              | -              | -              | -              |
| RippleNet | 0.1336         | 0.0910         | 0.0791         | 0.1238         | <b>0.0664</b>  | <b>0.0822</b>  |
| GC-MC     | 0.1316         | 0.0874         | 0.0818         | <b>0.1253</b>  | 0.0659         | 0.0790         |
| KGAT      | <b>0.1489*</b> | <b>0.1006*</b> | <b>0.0870*</b> | <b>0.1325*</b> | <b>0.0712*</b> | <b>0.0867*</b> |
| %Improv.  | 8.95%          | 10.05%         | 4.93%          | 5.77%          | 7.18%          | 5.54%          |

KGAT achieves the best performance

# Applications



## 1) Knowledge-enhanced Graph Learning for Recommendation

- Path-based recommendation methods
- Propagation-based recommendation methods




## 2) Knowledge-enhanced Graph Learning for Natural Language Processing (NLP)


- Natural Language Understanding
- Commonsense Reasoning
- Advancing LLMs with Knowledge

# Knowledge-enhanced Graph Learning for NLP



- A language model (LM) learns **how to express**

I go school to to want. 

I want to go to school. 

- Knowledge indicates **what to express**

Q: Where is the painting **Mona Lisa**?

A: It is in **Louvre, Paris**.





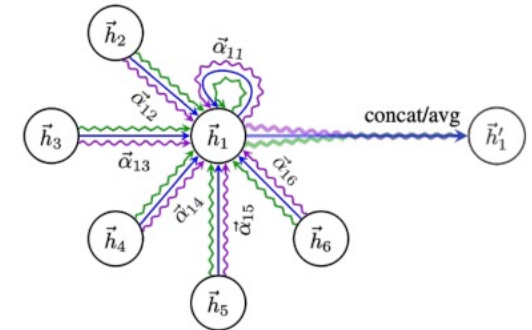
# Knowledge-enhanced Graph Learning for NLP



## To integrate knowledge into Language Models

- Step 1: Ground language into related knowledge
  - String matching, NER, Entity linking, information retrieval
  - Identify concepts and relations in the knowledge source
- Step 2: Represent knowledge
  - Concept descriptions, Graph embeddings
- Step 3: Fuse knowledge representation into language model
  - Concatenate concept descriptions into input
  - Append embeddings into input embeddings

The **pen** is on the **desk**.

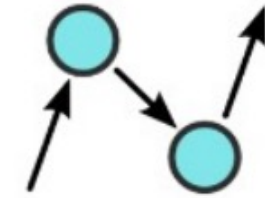


# Knowledge-enhanced Graph Learning for NLP

What are the knowledge sources for NLP



WIKIDATA



ConceptNet



WIKTIONARY  
the open content based dictionary



Domain-specific  
knowledge

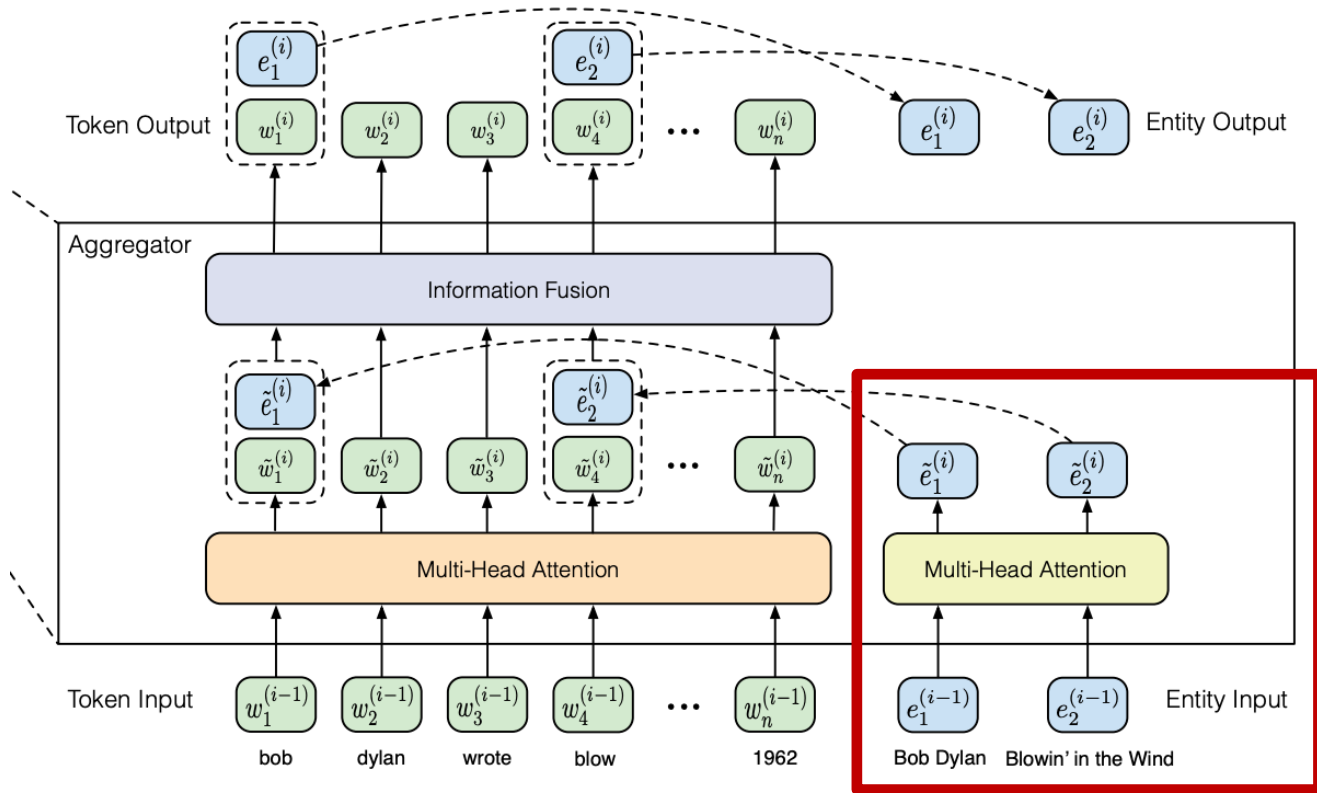


Commonsense  
knowledge

Wikipedia-based  
knowledge

# Knowledge-enhanced Graph Learning for NLP

## Natural Language Understanding: ERNIE



**Bob Dylan** wrote **Blowin' in the Wind** in 1962

- Use TransE to compute entity embeddings from Wikidata
- Concatenate token input and entity embeddings for learning

# Knowledge-enhanced Graph Learning for NLP



## Natural Language Understanding: ERNIE

| Model             | Acc.         | Macro        | Micro        |
|-------------------|--------------|--------------|--------------|
| NFGEC (Attentive) | 54.53        | 74.76        | 71.58        |
| NFGEC (LSTM)      | 55.60        | 75.15        | 71.73        |
| BERT              | 52.04        | 75.16        | 71.63        |
| ERNIE             | <b>57.19</b> | <b>76.51</b> | <b>73.39</b> |

Entity typing task on the FIGER dataset

| Model   | FewRel |       |              | TACRED |       |              |
|---------|--------|-------|--------------|--------|-------|--------------|
|         | P      | R     | F1           | P      | R     | F1           |
| CNN     | 69.51  | 69.64 | 69.35        | 70.30  | 54.20 | 61.20        |
| PA-LSTM | -      | -     | -            | 65.70  | 64.50 | 65.10        |
| C-GCN   | -      | -     | -            | 69.90  | 63.30 | 66.40        |
| BERT    | 85.05  | 85.11 | 84.89        | 67.23  | 64.81 | 66.00        |
| ERNIE   | 88.49  | 88.44 | <b>88.32</b> | 69.97  | 66.08 | <b>67.97</b> |

Relation classification task

| Model                | MNLI-(m/mm)<br>392k | QQP<br>363k | QNLI<br>104k | SST-2<br>67k |
|----------------------|---------------------|-------------|--------------|--------------|
| BERT <sub>BASE</sub> | 84.6/83.4           | 71.2        | -            | 93.5         |
| ERNIE                | 84.0/83.2           | 71.2        | 91.3         | 93.5         |

| Model                | CoLA<br>8.5k | STS-B<br>5.7k | MRPC<br>3.5k | RTE<br>2.5k |
|----------------------|--------------|---------------|--------------|-------------|
| BERT <sub>BASE</sub> | 52.1         | 85.8          | 88.9         | 66.4        |
| ERNIE                | 52.3         | 83.2          | 88.2         | 68.8        |

General Language Understanding Evaluation (GLUE) benchmark

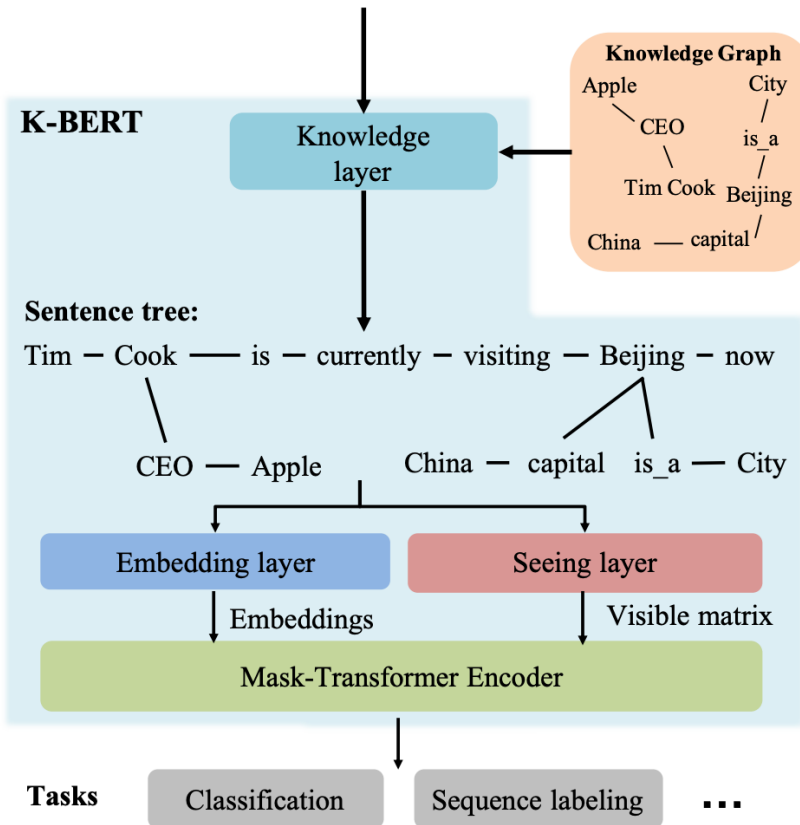
ERNIE outperforms BERT across tasks

# Knowledge-enhanced Graph Learning for NLP



## Natural Language Understanding: K-BERT

**Input sentence:** Tim Cook is currently visiting Beijing now



- K-BERT injects relevant triples from the KG and transform the original sentence into a knowledge-rich sentence tree for encoding
- K-BERT is equipped with an editable KG, which can be adapted to its application domain. For example, a medical KG can be used to grant the K-BERT with medical knowledge

# Knowledge-enhanced Graph Learning for NLP

## Natural Language Understanding: K-BERT



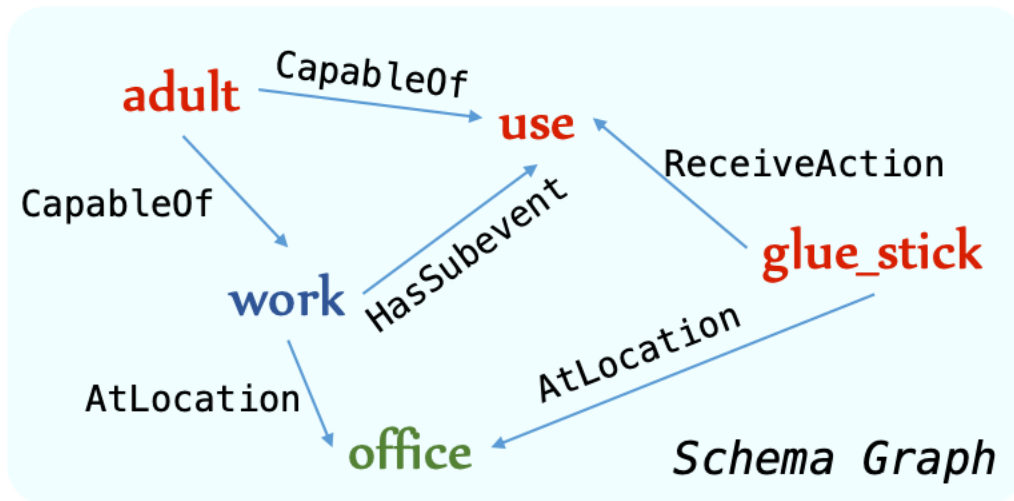
Results on different domain-specific tasks

| Models\Datasets                  | Finance_Q&A |           |             | Law_Q&A   |           |             | Finance_NER |           |             | Medicine_NER |           |             |
|----------------------------------|-------------|-----------|-------------|-----------|-----------|-------------|-------------|-----------|-------------|--------------|-----------|-------------|
|                                  | <i>P.</i>   | <i>R.</i> | <i>F1</i>   | <i>P.</i> | <i>R.</i> | <i>F1</i>   | <i>P.</i>   | <i>R.</i> | <i>F1</i>   | <i>P.</i>    | <i>R.</i> | <i>F1</i>   |
| Pre-trained on WikiZh by Google. |             |           |             |           |           |             |             |           |             |              |           |             |
| Google BERT                      | 81.9        | 86.0      | 83.9        | 83.1      | 90.1      | 86.4        | 84.8        | 87.4      | 86.1        | 91.9         | 93.1      | 92.5        |
| K-BERT (HowNet)                  | 83.3        | 84.4      | 83.9        | 83.7      | 91.2      | 87.3        | 86.3        | 89.0      | <b>87.6</b> | 93.2         | 93.3      | 93.3        |
| K-BERT (CN-DBpedia)              | 81.5        | 88.6      | <b>84.9</b> | 82.1      | 93.8      | <b>87.5</b> | 86.1        | 88.7      | 87.4        | 93.9         | 93.8      | 93.8        |
| K-BERT (MedicalKG)               | -           | -         | -           | -         | -         | -           | -           | -         | -           | 94.0         | 94.4      | <b>94.2</b> |

K-BERT benefits from the domain knowledge in KGs and performs well across different domain-specific tasks

# Knowledge-enhanced Graph Learning for NLP

## Commonsense reasoning: KagNet



Grounding ↑ Knowledge-Aware Commonsense Inference ↓

Where do adults use glue sticks?

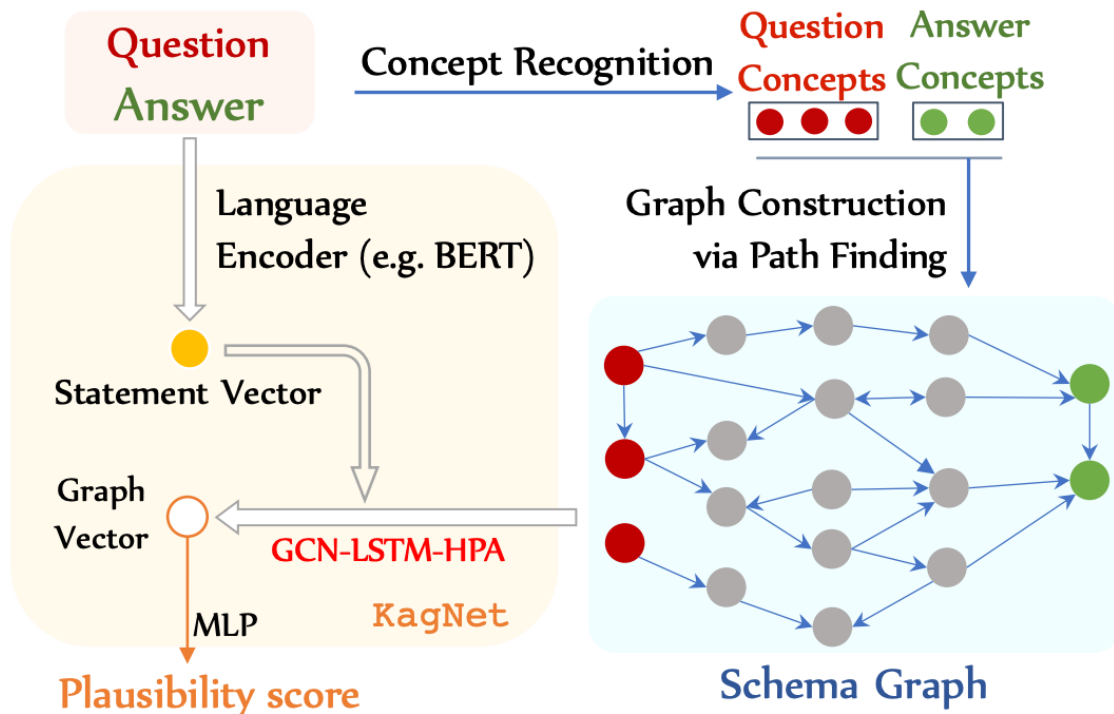
A: classroom B: office C: desk drawer

An example of using external commonsense knowledge for inference in natural language commonsense questions



# Knowledge-enhanced Graph Learning for NLP

## Commonsense reasoning: KagNet



- For each pair of question and answer candidate, KagNet retrieves a graph from external knowledge graphs (e.g. ConceptNet)
- The retrieved graph contains relevant knowledge for determining the plausibility of a given answer choice

# Knowledge-enhanced Graph Learning for NLP

## Commonsense reasoning: KagNet



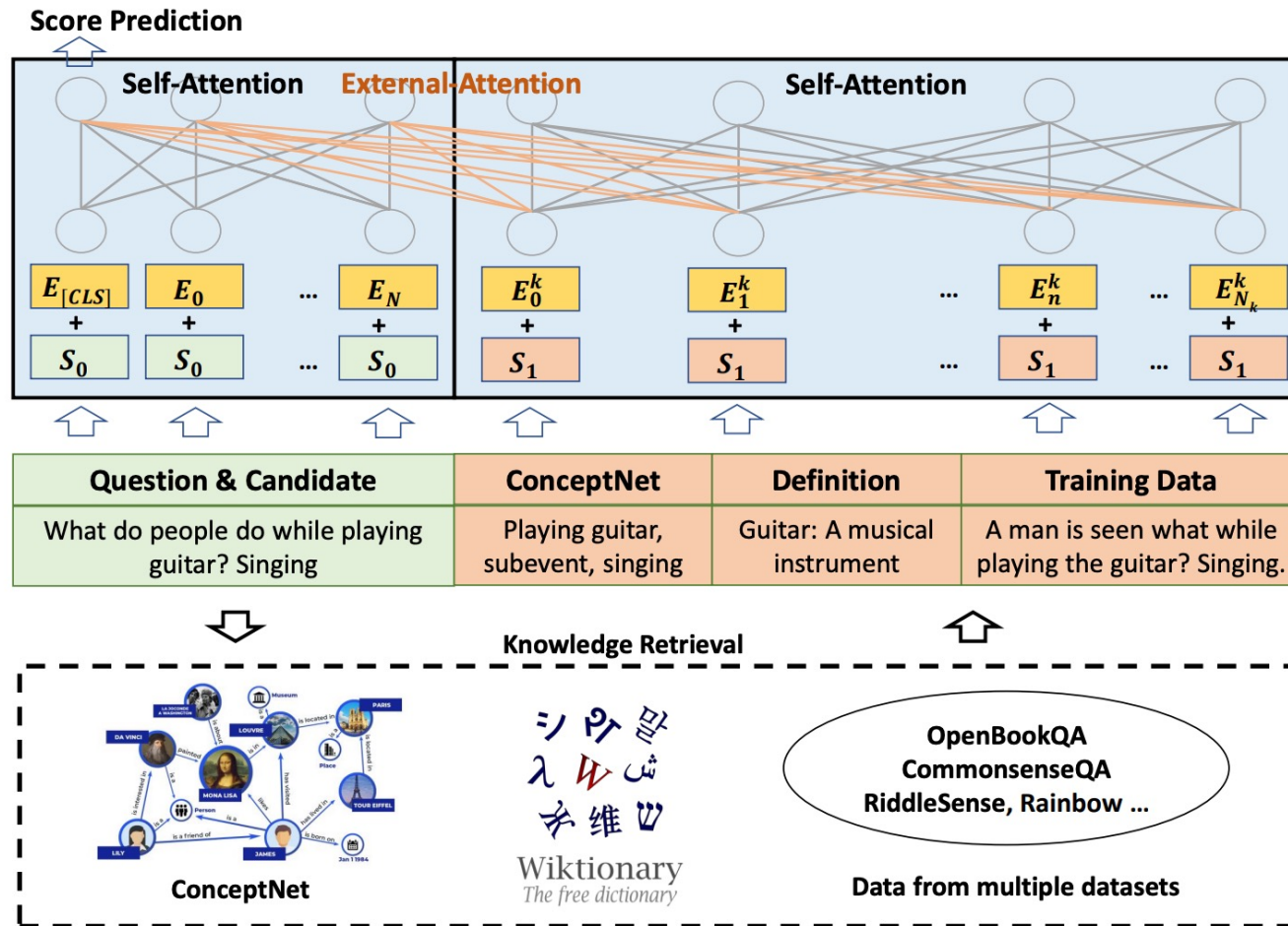
| Model                 | 10(%) of IHtrain |                | 50(%) of IHtrain |                | 100(%) of IHtrain |                |
|-----------------------|------------------|----------------|------------------|----------------|-------------------|----------------|
|                       | IHdev-Acc.(%)    | IHtest-Acc.(%) | IHdev-Acc.(%)    | IHtest-Acc.(%) | IHdev-Acc.(%)     | IHtest-Acc.(%) |
| Random guess          | 20.0             | 20.0           | 20.0             | 20.0           | 20.0              | 20.0           |
| GPT-FINETUNING        | 27.55            | 26.51          | 32.46            | 31.28          | 47.35             | 45.58          |
| GPT-KAGNET            | 28.13            | <b>26.98</b>   | 33.72            | <b>32.33</b>   | 48.95             | <b>46.79</b>   |
| BERT-BASE-FINETUNING  | 30.11            | 29.78          | 38.66            | 36.83          | 53.48             | 53.26          |
| BERT-BASE-KAGNET      | 31.05            | <b>30.94</b>   | 40.32            | <b>39.01</b>   | 55.57             | <b>56.19</b>   |
| BERT-LARGE-FINETUNING | 35.71            | 32.88          | 55.45            | 49.88          | 60.61             | 55.84          |
| BERT-LARGE-KAGNET     | 36.82            | <b>33.91</b>   | 58.73            | <b>51.13</b>   | 62.35             | <b>57.16</b>   |
| Human Performance     | -                | 88.9           | -                | 88.9           | -                 | 88.9           |

Using KagNet outperforms fine-tuning the language models themselves

# Knowledge-enhanced Graph Learning for NLP



## Commonsense reasoning: KEAR



Related knowledge is retrieved from external sources, e.g., **knowledge graph, dictionary and training data**, using the input as the key and then integrated with the input

# Knowledge-enhanced Graph Learning for NLP

## Commonsense reasoning: KEAR



| Method          | E-I+VAT | D-xxl | DV3-I |
|-----------------|---------|-------|-------|
| Base            | 82.1    | 83.8  | 84.6  |
| + KG            | 85.2    | 86.4  | 86.7  |
| + Dictionary    | 83.8    | 84.0  | 85.1  |
| + Training data | 84.0    | 86.4  | 87.1  |

**E-I+VAT:** ELECTRA-large with VAT

**D-xxl:** DeBERTa-xxlarge

**DV3-I:** DeBERTaV3-large

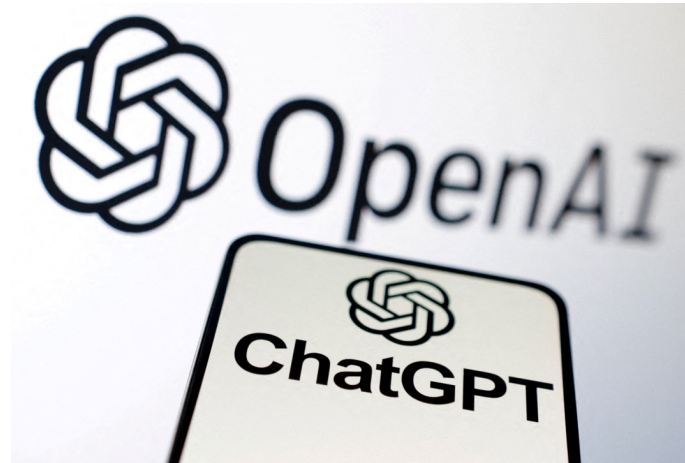
Results of applying external attention to different knowledge sources

- These knowledge sources bring gains in commonsense reasoning accuracy across all base encoder models
- Generally, knowledge graphs contribute the most

# Advancing LLMs with Knowledge

## What are LLMs (Large Language Models)?

A type of machine learning model that can perform a variety of language tasks.



# Advancing LLMs with Knowledge



## Why LLMs Need External Knowledge?

⚠ Limitations

The problem of hallucination →

May occasionally generate incorrect information

(from OpenAI)

YI

**You**

Were there fossil fuels in the ground when humans evolved?

~300 Million Years Ago

~300 Thousand Years Ago

Options:

- (a) This was only created by humans
- (b) Humans predate fossil fuel formation
- (c) Significant supplies accumulated prior
- (d) None of these

→ **Correct Answer**



# Advancing LLMs with Knowledge

## Why LLMs Need External Knowledge?



### The problem of hallucination



You

Were there fossil fuels in the ground when humans evolved?

Options:

- (a) This was only created by humans
- (b) Humans predate fossil fuel formation
- (c) Significant supplies accumulated prior
- (d) None of these



ChatGPT

(b) Humans predate fossil fuel formation

Wrong Answer ←

Explanations that  
Make No Sense ←

Fossil fuels, such as coal, oil, and natural gas, are formed from the remains of ancient plants and microorganisms that lived millions of years ago. The process of fossil fuel formation takes a very long time, and it began long after the evolution of early humans. Therefore, humans predate the formation of fossil fuels, making option (b) the correct answer.





# Advancing LLMs with Knowledge

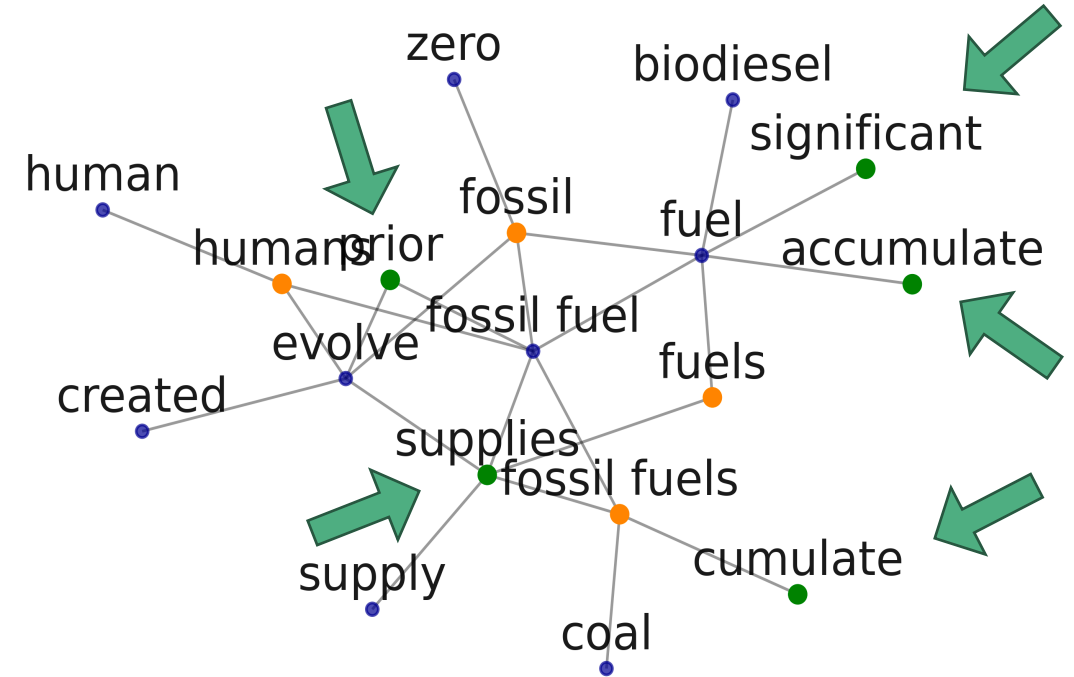
## How External Knowledge Assists?



Provides Factual Information

Were there **fossil fuels** in the ground when **humans** evolved?

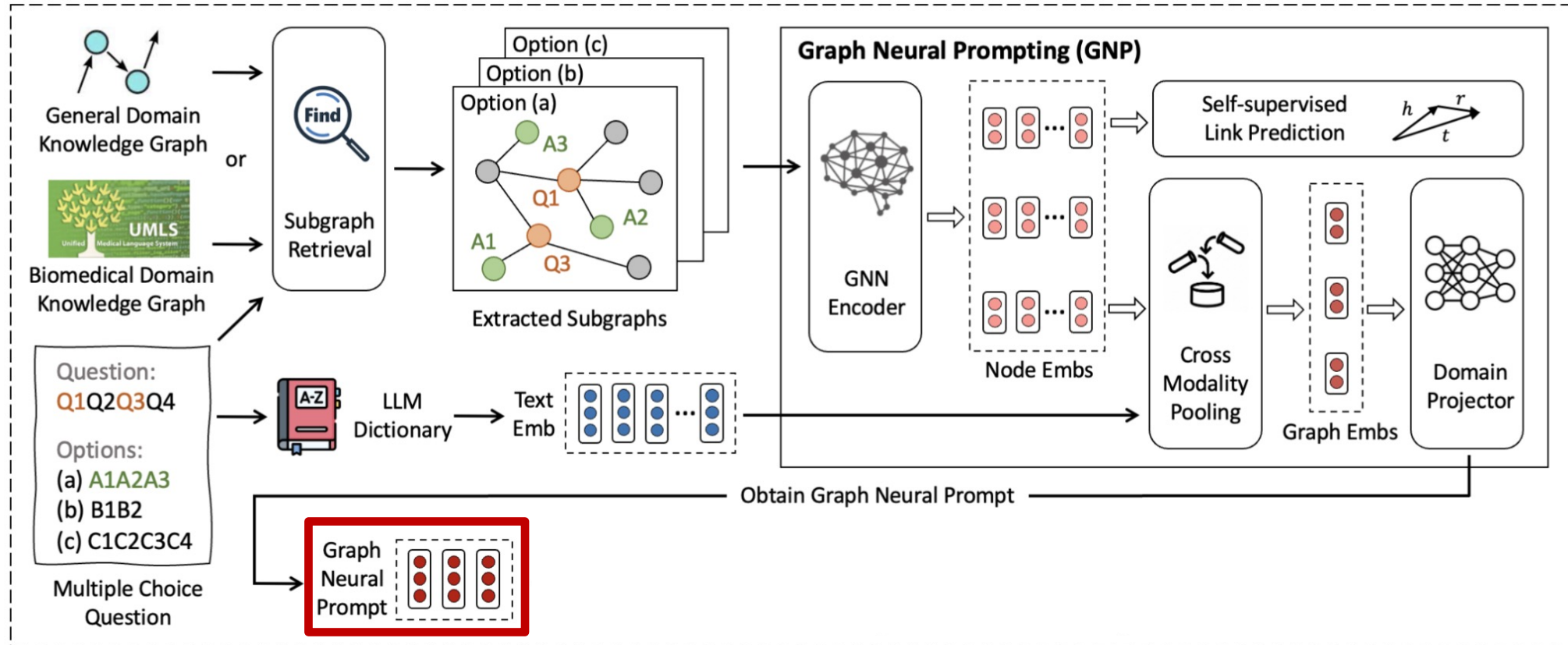
(c) Significant supplies accumulated prior



Retrieved Knowledge for Question  
from ConceptNet

# Advancing LLMs with Knowledge

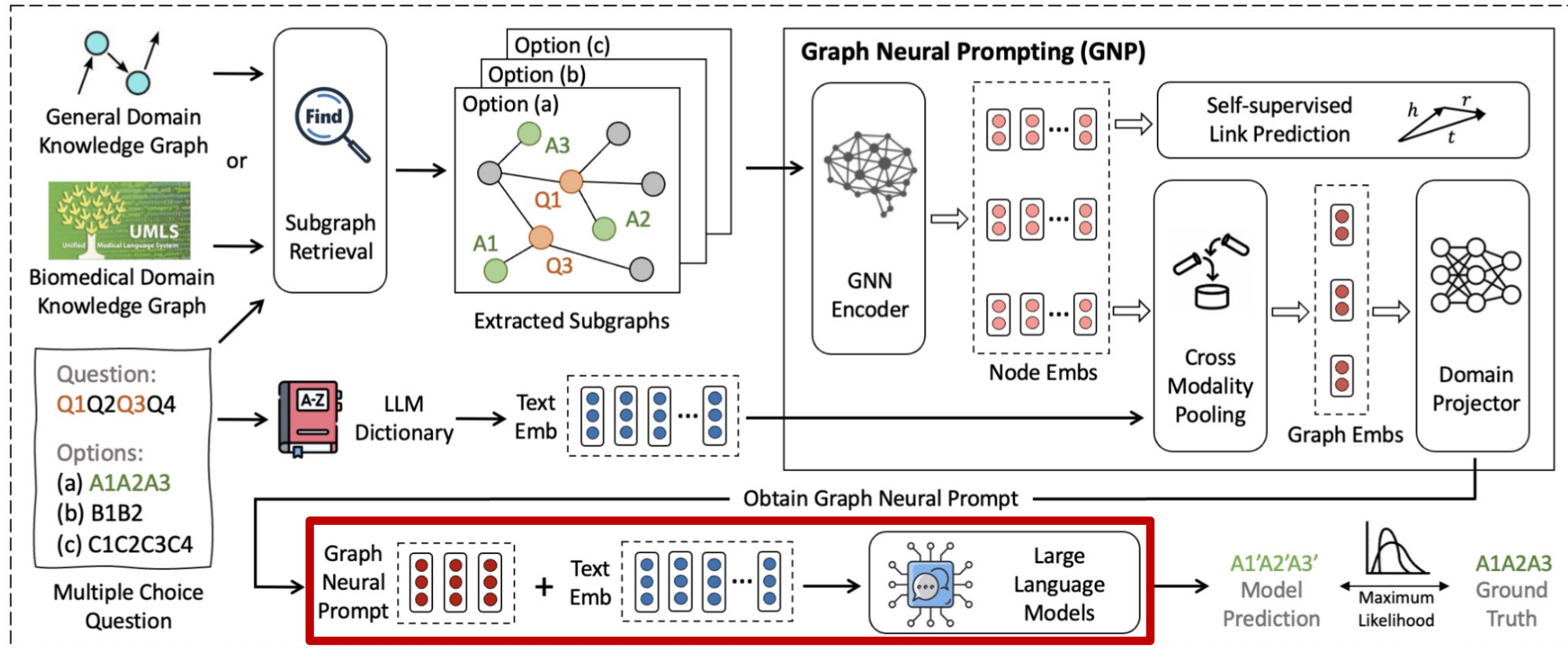
## Graph Neural Prompting (GNP)



Extracted Knowledge

# Advancing LLMs with Knowledge

## Graph Neural Prompting (GNP)



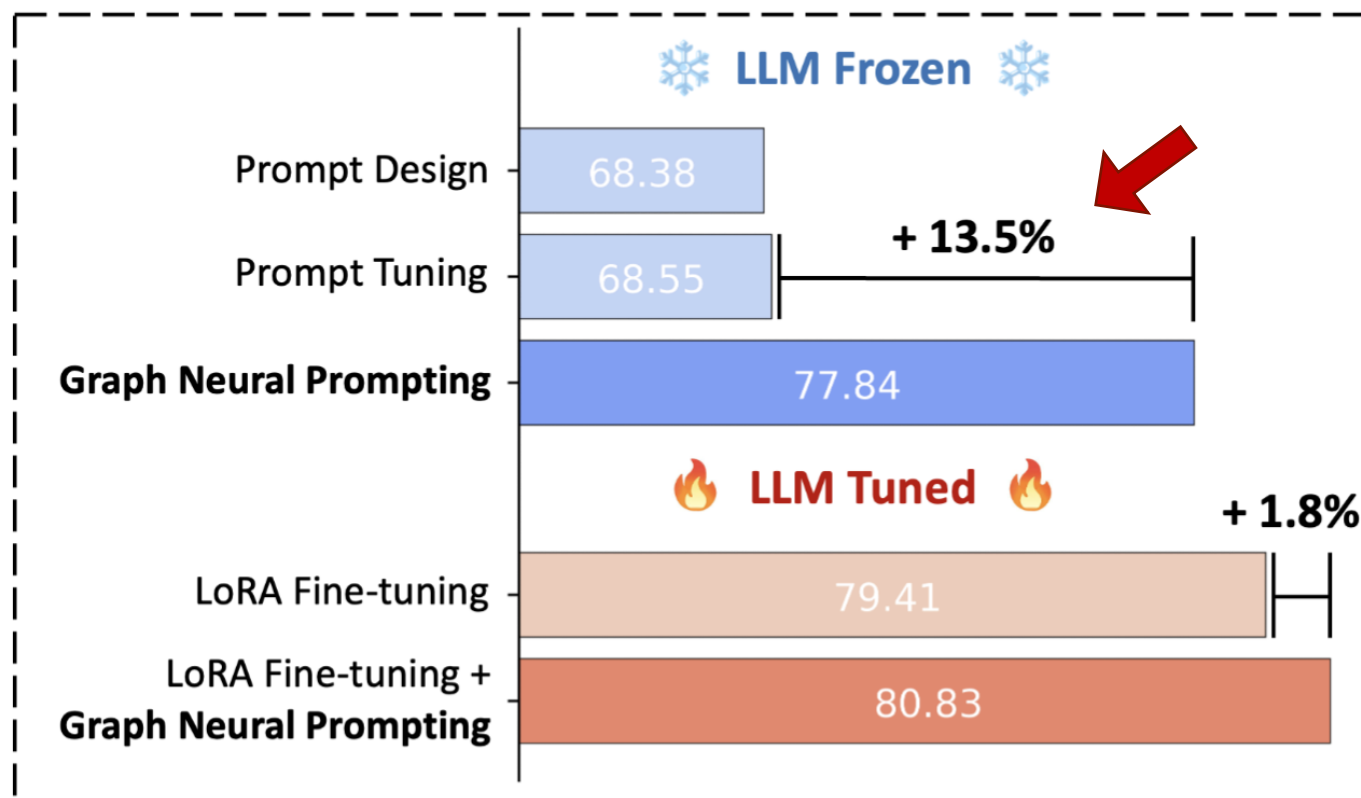
Extracted Knowledge → Real-life Applications

# Advancing LLMs with Knowledge

## Graph Neural Prompting (GNP)



w/o Knowledge  
w/ Knowledge



Significant Improvement

Results Averaged Across 4 Commonsense and 2 Biomedical Datasets

# Applications

## Chapter summary



### 1) Knowledge-enhanced Graph Learning for Recommendation

- Path-based recommendation methods
- Propagation-based recommendation methods

### 2) Knowledge-enhanced Graph Learning for Natural Language Processing (NLP)

- Natural Language Understanding
- Commonsense Reasoning
- Advancing LLMs with Knowledge

# Tutorial Outline



- Preliminaries and Foundations
- Graph Learning Enhanced by Knowledge from Data
- Graph Learning Enhanced by Knowledge from Models
- Graph Learning Enhanced by Knowledge from Humans and Domains
- Graph Learning Enhanced by Knowledge from External Sources
- Knowledge-enhanced Graph Learning for Real-world Applications

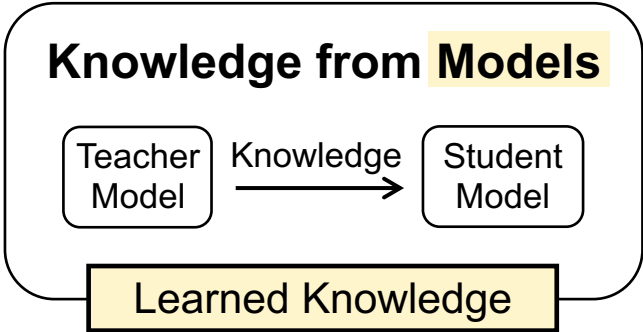
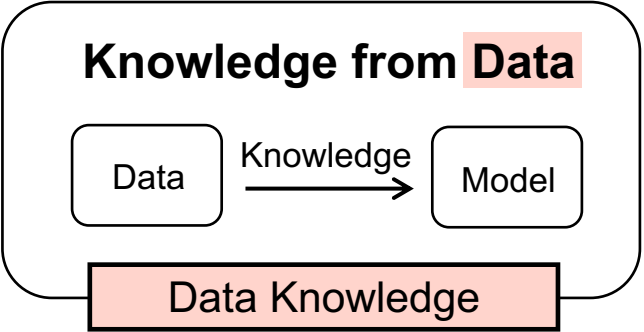


- **Summary and Future Directions**

# Summary: Knowledge-enhanced Graph Learning

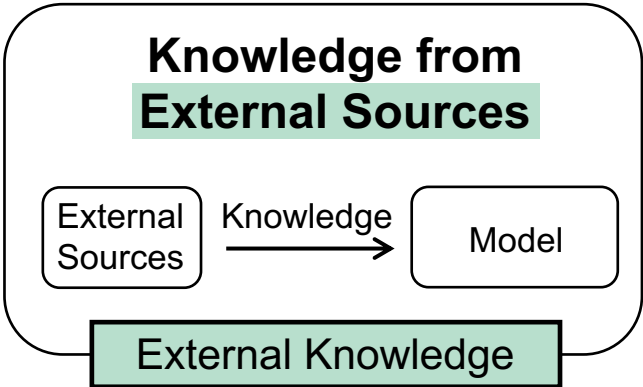
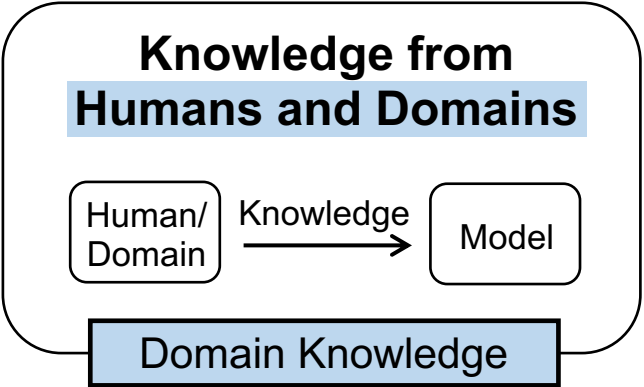


Implicit Knowledge



Real-world Applications

Explicit Knowledge



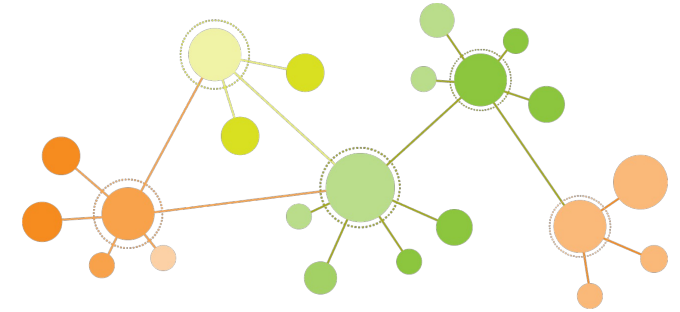
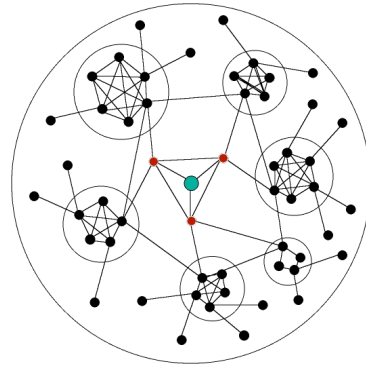
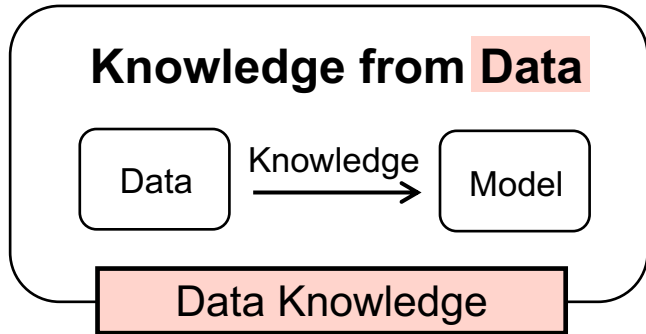
Future Directions



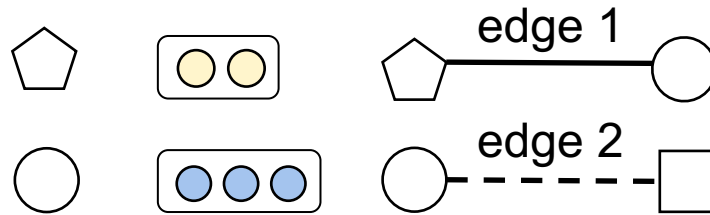
# Summary: Knowledge-enhanced Graph Learning



## Knowledge from data



## Local communities



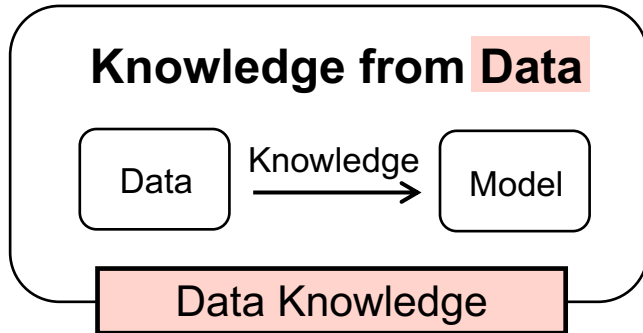
## Heterogeneous Information

## Node positions

|    | p1 | p2 | p3 |           | p1 | p2 | p3 |   |
|----|----|----|----|-----------|----|----|----|---|
| p1 | 1  | 1  | 0  | Co-author | p1 | 1  | 1  | 0 |
| p2 | 1  | 1  | 1  |           | p2 | 1  | 1  | 0 |
| p3 | 0  | 1  | 1  |           | p3 | 0  | 0  | 1 |

## Higher-order Semantics

# Summary: Knowledge-enhanced Graph Learning



Knowledge can be obtained from

## 1) single-instance level perception

- **Node sampling** for node positions ([P-GNN](#))

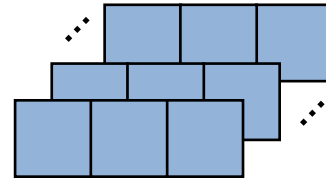
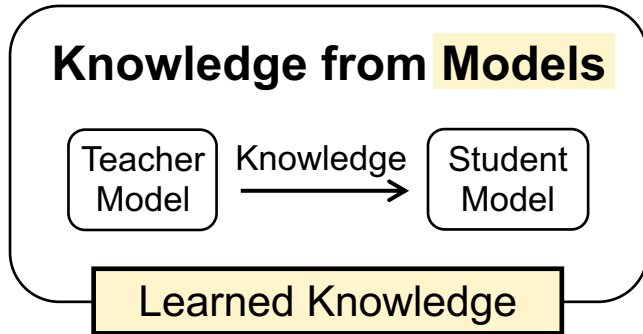
## 2) multiple-instance level perception

- **Path sampling** for positional and semantic information ([HGMAE](#))
- **Subgraph sampling** for community information ([SEAL](#), [WalkPool](#))

# Summary: Knowledge-enhanced Graph Learning



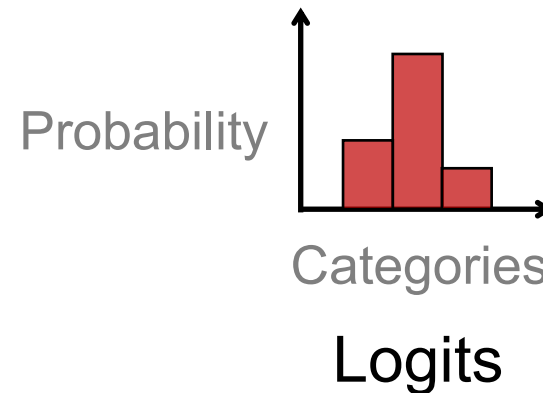
## Knowledge from models



Learned Embeddings

|    | p1 | p2 | p3 |
|----|----|----|----|
| p1 | 1  | .3 | .7 |
| p2 | .3 | 1  | .2 |
| p3 | .7 | .2 | 1  |

Data Similarity



# Summary: Knowledge-enhanced Graph Learning

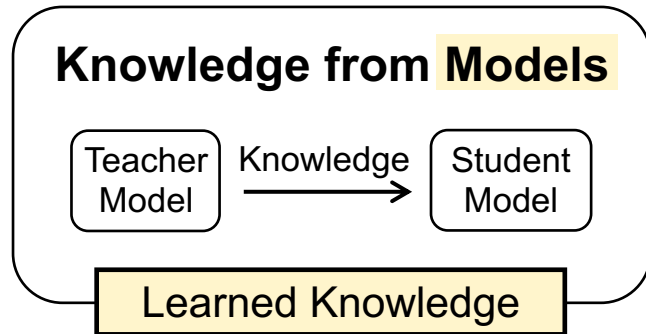


## How to obtain knowledge from

- Logits ([TinyGNN](#), [BGNN](#))
- Embeddings ([GraphAKD](#))
- Structures ([LSP](#), [G-CRD](#))

## A recent learning strategy

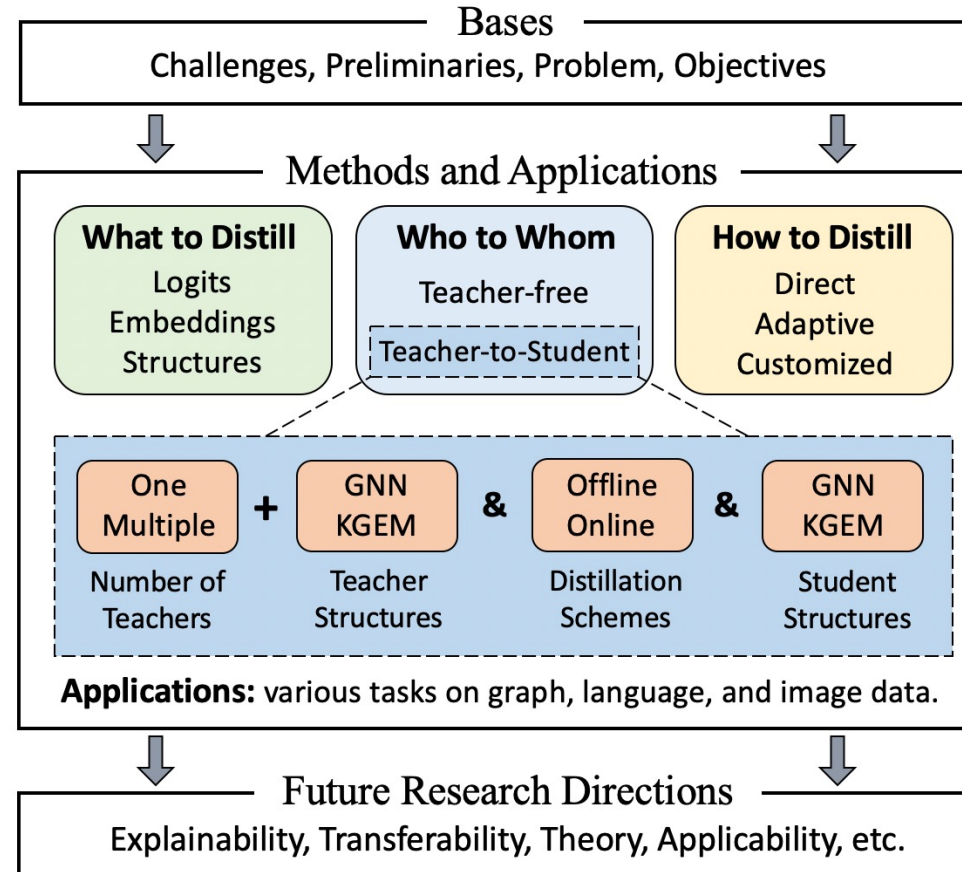
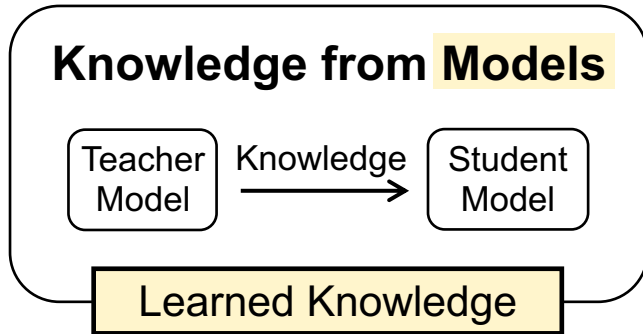
- distilling an MLP to replace GNNs ([GLNN](#), [NOSMOG](#))



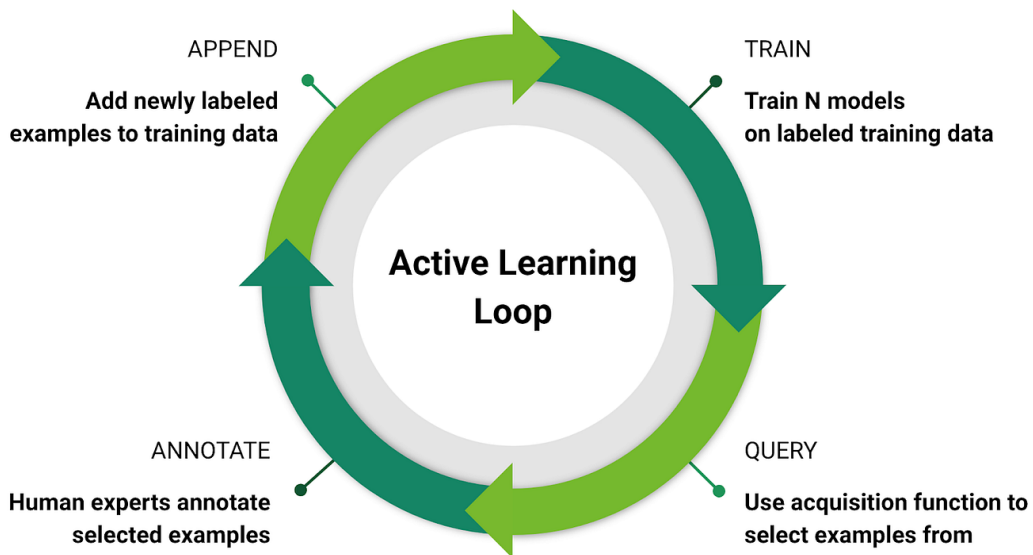
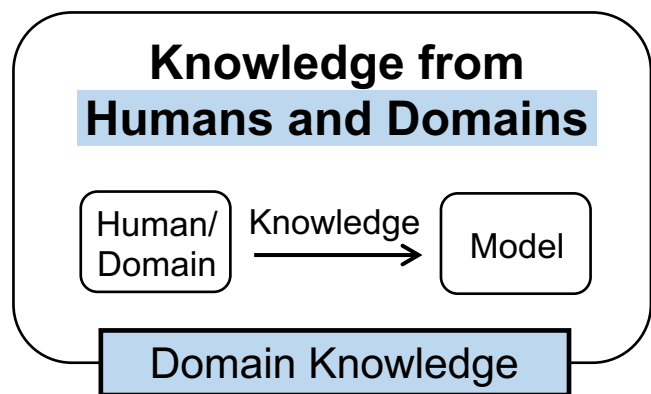
# Summary: Knowledge-enhanced Graph Learning



More details can be found in this survey



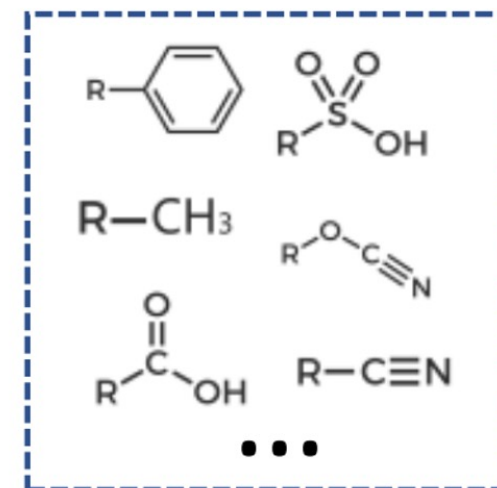
# Summary: Knowledge-enhanced Graph Learning



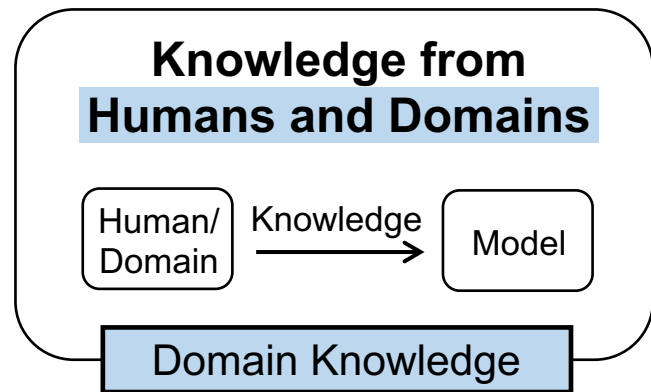
Graph active learning with human knowledge

Domain knowledge such as using motif to represent functional groups in chemistry

Semantic motifs from domain knowledge



# Summary: Knowledge-enhanced Graph Learning



## Graph learning enhanced by human feedback

- Graph active learning with human knowledge ([ALG](#), [IGP](#), [OWGAL](#))

## Graph learning enhanced by domain knowledge

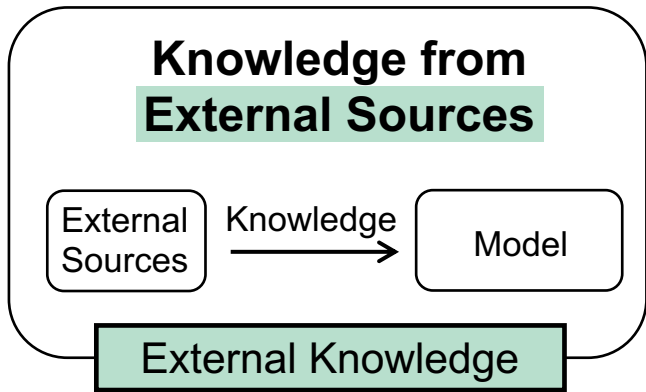
- Chemistry domain knowledge for molecular property prediction ([GROVER](#), [MGSSL](#))



# Summary: Knowledge-enhanced Graph Learning



## Knowledge from external sources



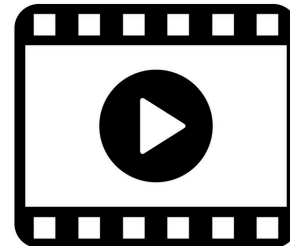
Text



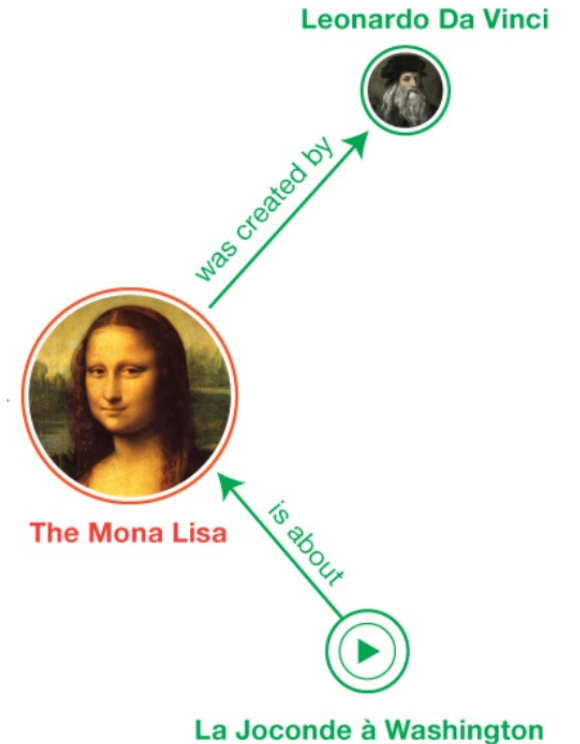
Wikipedia



Image

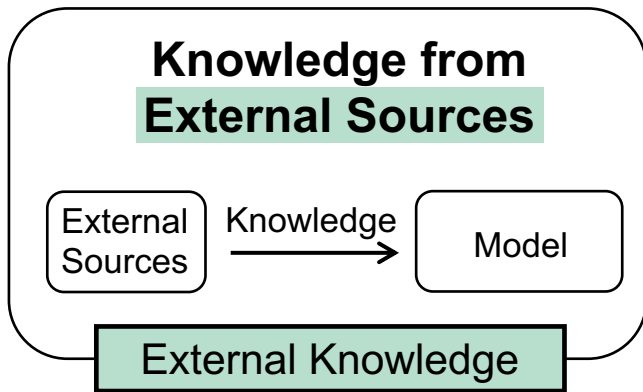


Video



Knowledge triplets

# Summary: Knowledge-enhanced Graph Learning



## Graph learning on text-rich graphs

- Graph learning on textual-node graphs ([PATTON](#), [Heterformer](#))
- Graph learning on textual-edge graphs ([Edgeformers](#))


## Graph learning on knowledge graphs

- Knowledge Graph Embedding ([TransE](#), [DistMult](#), [CompLex](#), [RotatE](#))
- Advancing KG tasks with text data ([KEPLER](#))

# Summary: Knowledge-enhanced Graph Learning



What are the benefits of utilizing knowledge

 Reduce reliance on massive data and intricate model

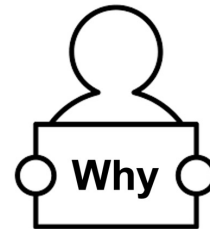
 Trustworthiness 

 Efficiency 

 Performance 



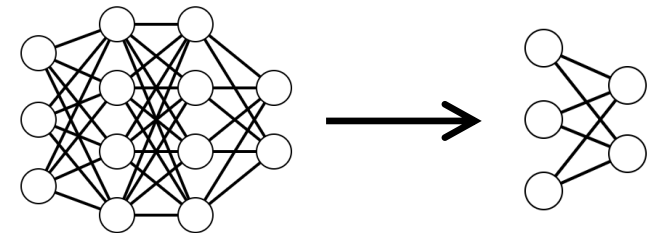
Robustness



Explainability

**NO LABELS**

Less labeled data



Smaller/simpler model

# Future Directions



## 1) How to obtain high-quality knowledge?

Knowledge sources can contain several issues such as Noises, biases, imbalances, missing or incorrect knowledge

Therefore, several research questions arise:

- How to detect and mitigate these issues?
- How to extract high-quality knowledge given these issues?
- How to verify the accuracy and authenticity of the knowledge?
- How to construct knowledge bases without these issues?

# Future Directions



## 2) How to Integrating knowledge from different sources?

Existing works mainly consider leveraging knowledge from one source, while ignoring the fact that various sources can provide different knowledge

Therefore, several research questions arise:

- How to identify the right knowledge sources?
- How to prioritize and extract knowledge from different sources?
- How to combine knowledge with different modalities or formats?
- How to integrate knowledge in a complementary manner?

# Future Directions



## 3) How to handle new knowledge?

New knowledge emerges incrementally and continuously in different time periods

Therefore, several research questions arise:

- How to collect, measure, and evaluate new knowledge?
- How to dynamically update new knowledge into knowledge base?
- How to forget old and irrelevant knowledge?
- How to balance and encode new knowledge?



AAAI-24 Tutorial  
Feb 2024, Vancouver



Thanks for listening!

Tutorial Website:



[yijuntian.com/tutorial](https://yijuntian.com/tutorial)

Yijun is looking for  
job opportunities