

Semantic-Structural Integration in Text-Attributed Graphs

Yuan FANG

School of Computing and Information Systems
Singapore Management University

ISCSC 2025 @ Taiyuan, China
9 Nov 2025

Outline

2

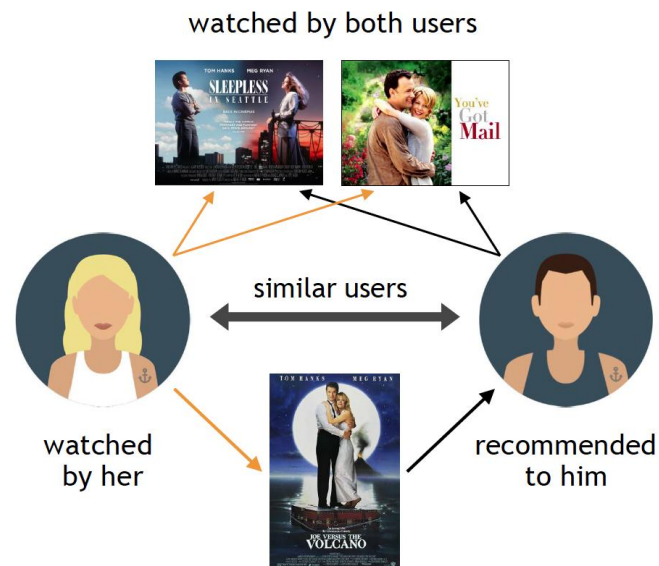
- **Introduction: Graphs & text-attributed graphs**
- Jointly training graph and textual data
- Quantizing graphs into language tokens for LLMs
- Conclusions

Graph structures are prevalent

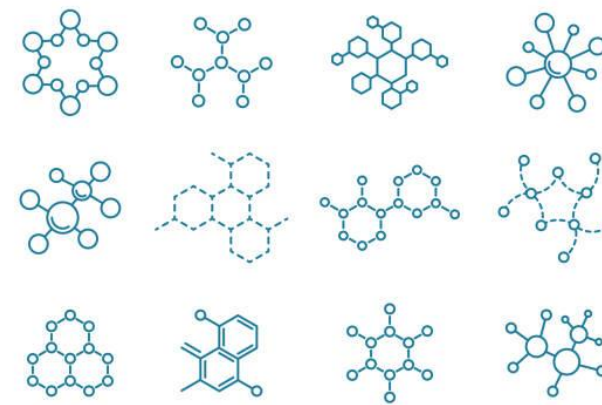
3



Social network



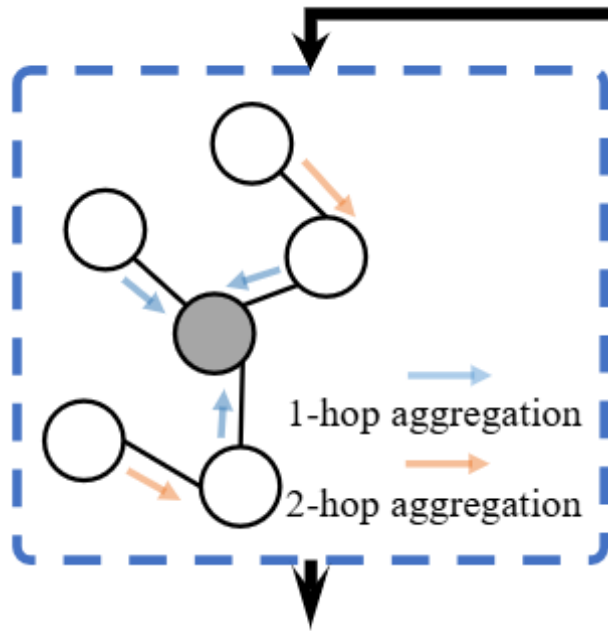
Recommendation System



Molecular graphs

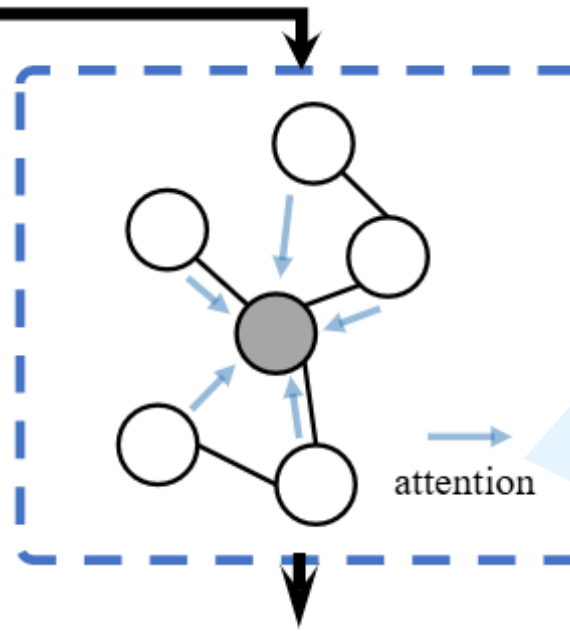
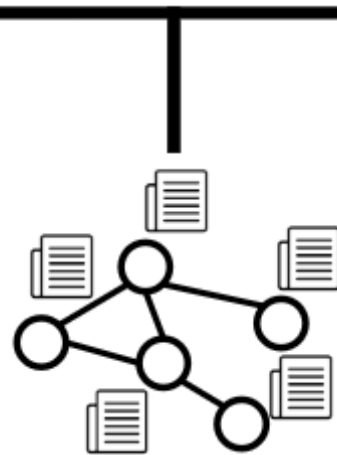
Learning from graph structures

4



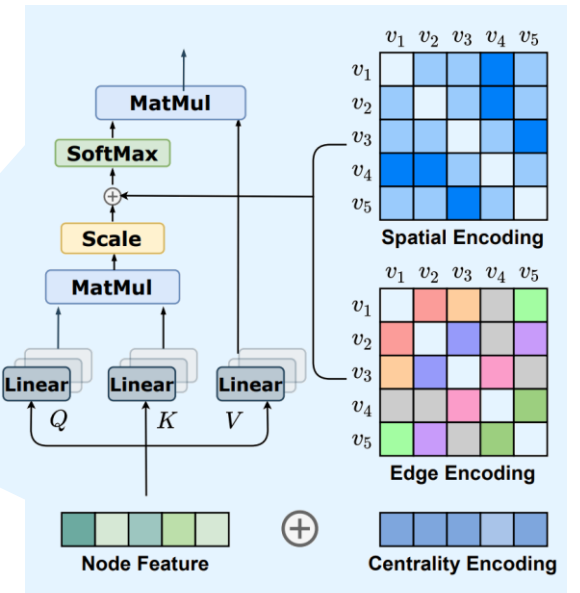
Predictions

**Message-passing
graph neural networks (GNNs)**



Predictions

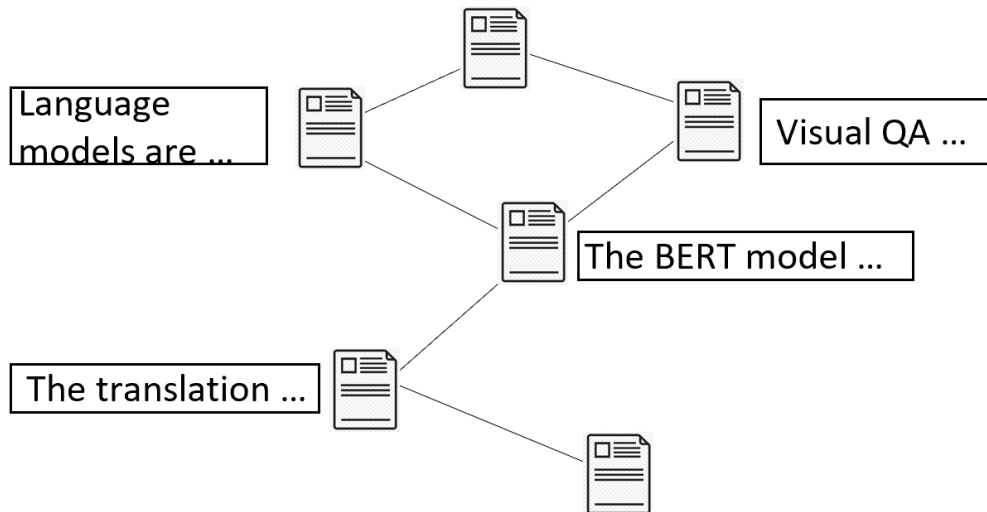
Graph transformers



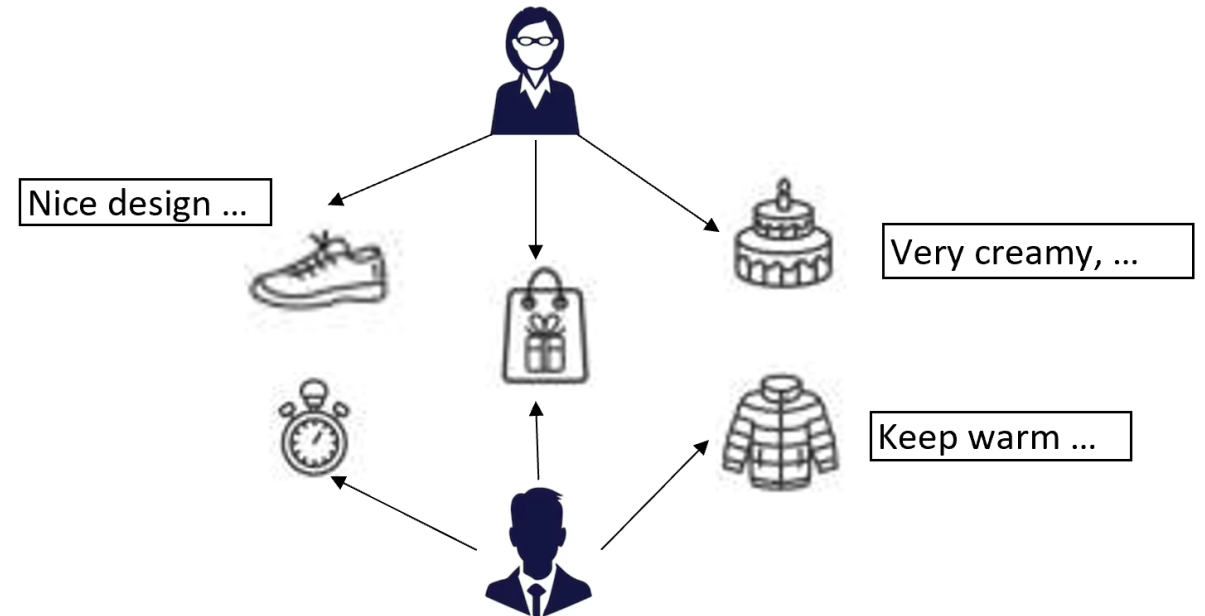
Ying, *et al.* Do Transformers Really Perform Bad for Graph Representation? NeurIPS 2021.
Liu, *et al.* Graph foundation models: Concepts, Opportunities and Challenges. TPAMI 2025.

Semantics on graphs: Text-Attributed Graphs

5



Citation graph for online articles

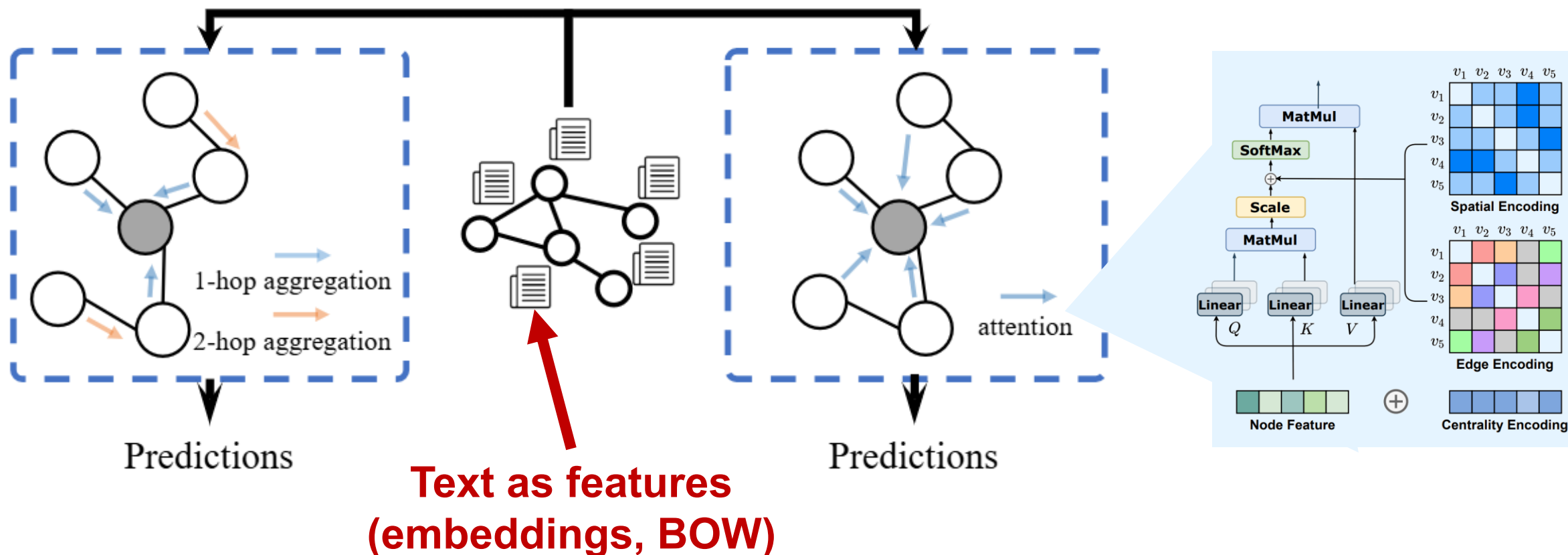


E-commerce item review graph

Can we **integrate** graph structures and textual semantics within one model?

Can GNNs or graph transformers utilize textual attributes? Yes, but ineffective

6



Ying, *et al.* Do Transformers Really Perform Bad for Graph Representation? NeurIPS 2021.

Liu, *et al.* Graph foundation models: Concepts, Opportunities and Challenges. TPAMI 2025.

Outline

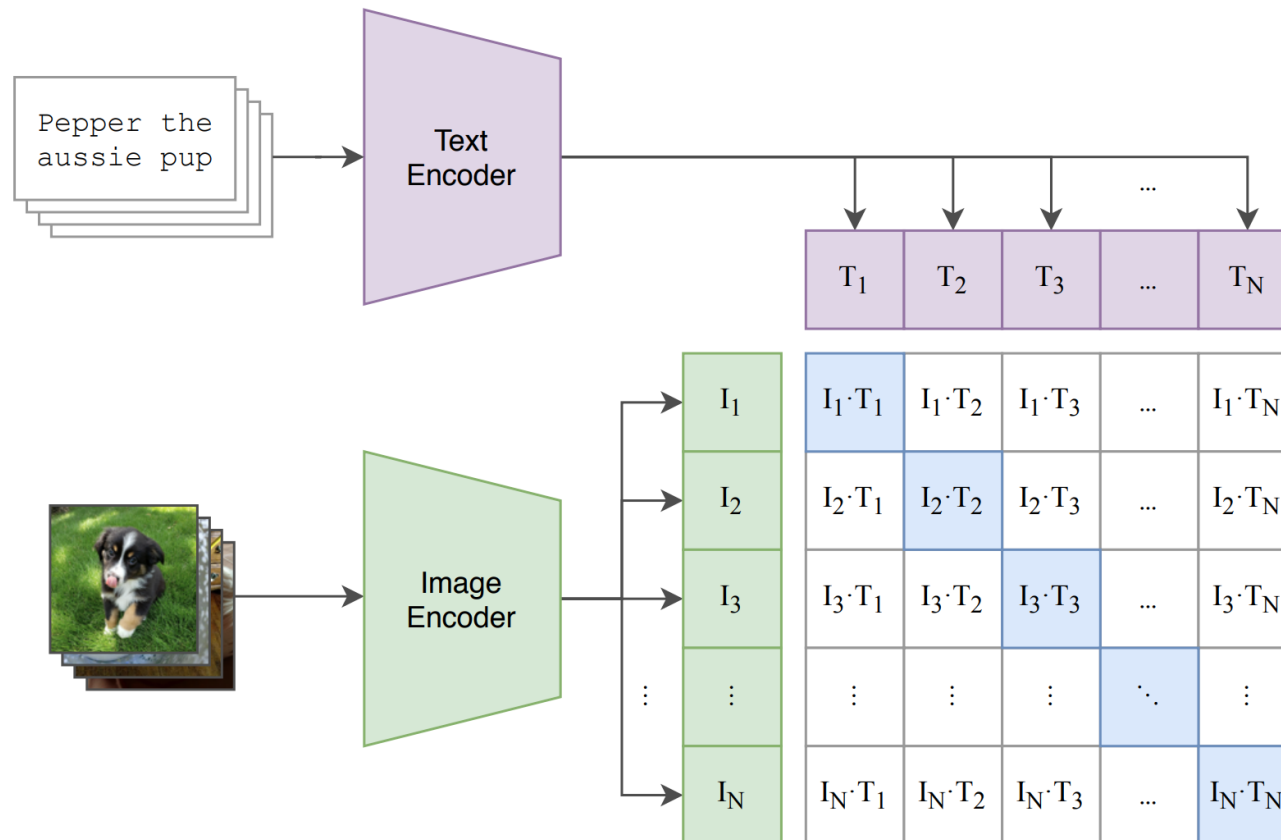
7

- Introduction: Graphs & text-attributed graphs
- **Jointly training graph and textual data**
- Quantizing graphs into language tokens for LLMs
- Conclusions

How are language-image models trained?

8

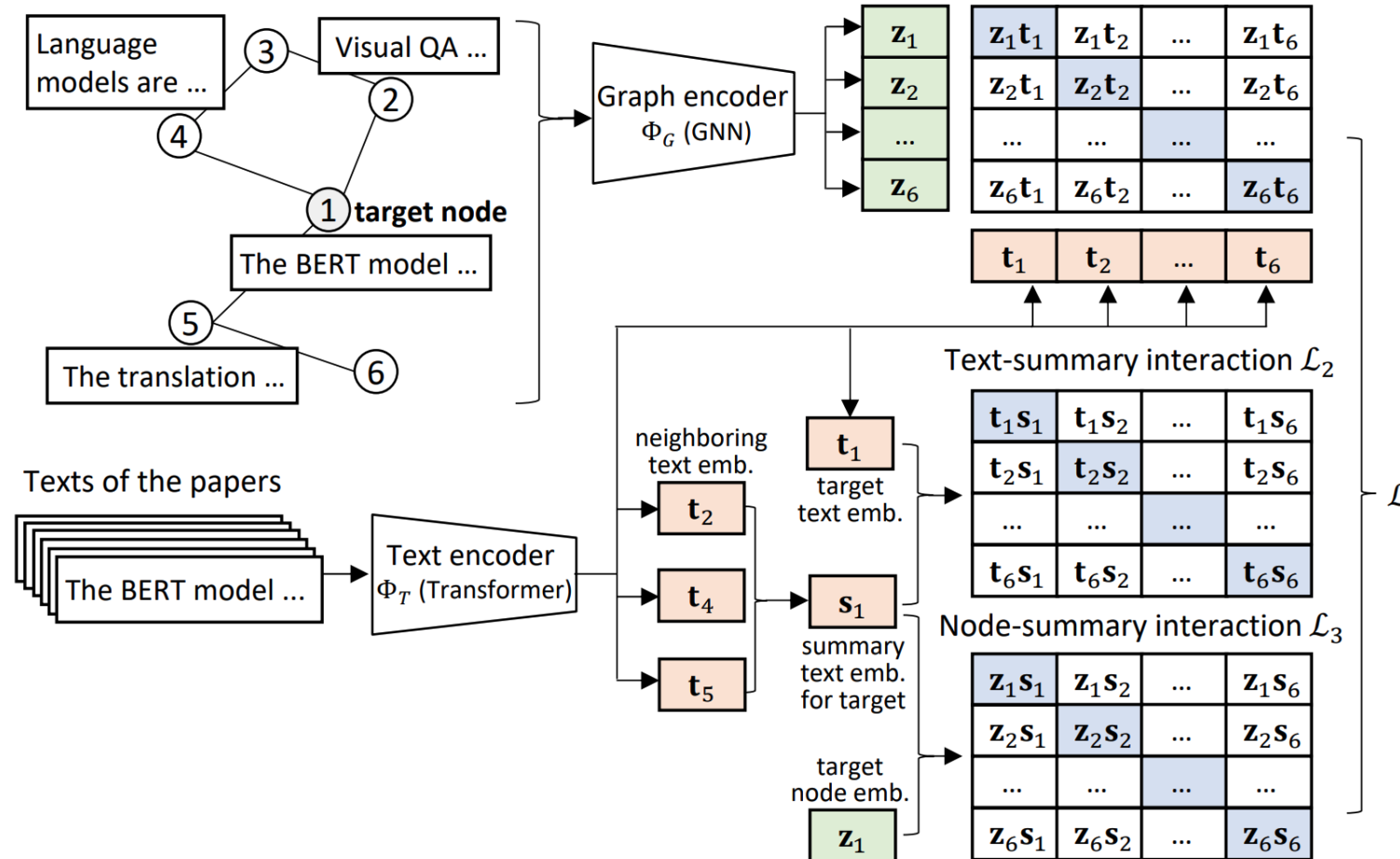
CLIP: Contrastive Language-Image Pre-training



Graph-grounded pre-training and prompting (G2P2)

9

Papers grounded on a citation network



Learns a dual-modal embedding space by jointly pre-training a **text encoder** and **graph encoder**

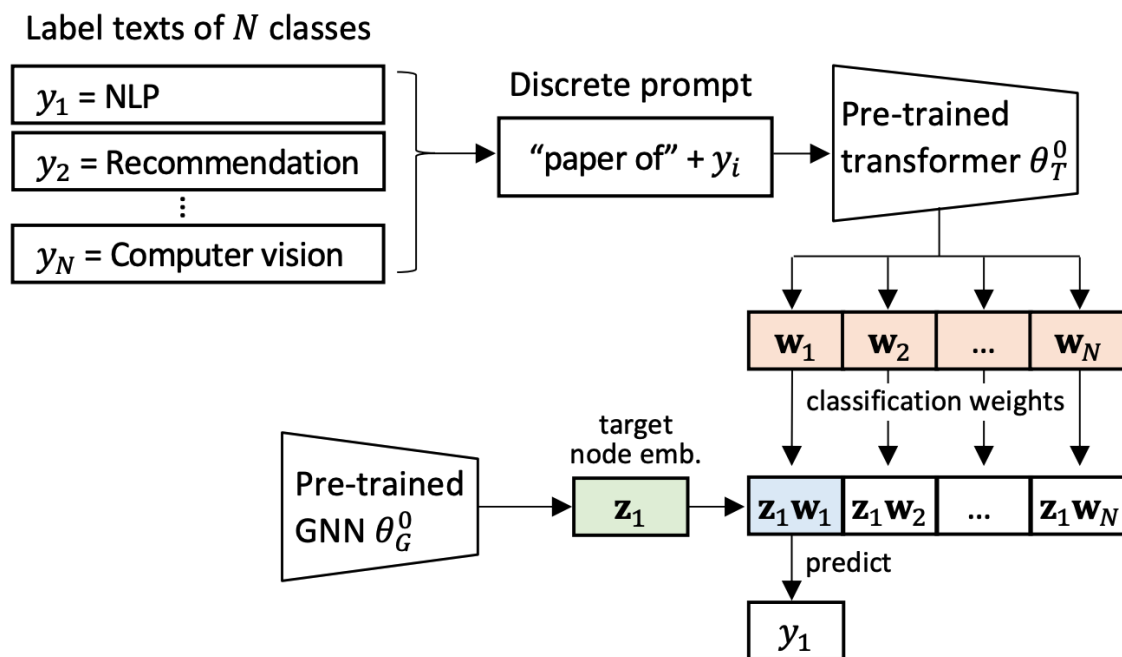
Exploits **three contrastive losses**

- \mathcal{L}_1 : Text-node contrast
- \mathcal{L}_2 : Text-summary contrast
- \mathcal{L}_3 : Node-summary contrast

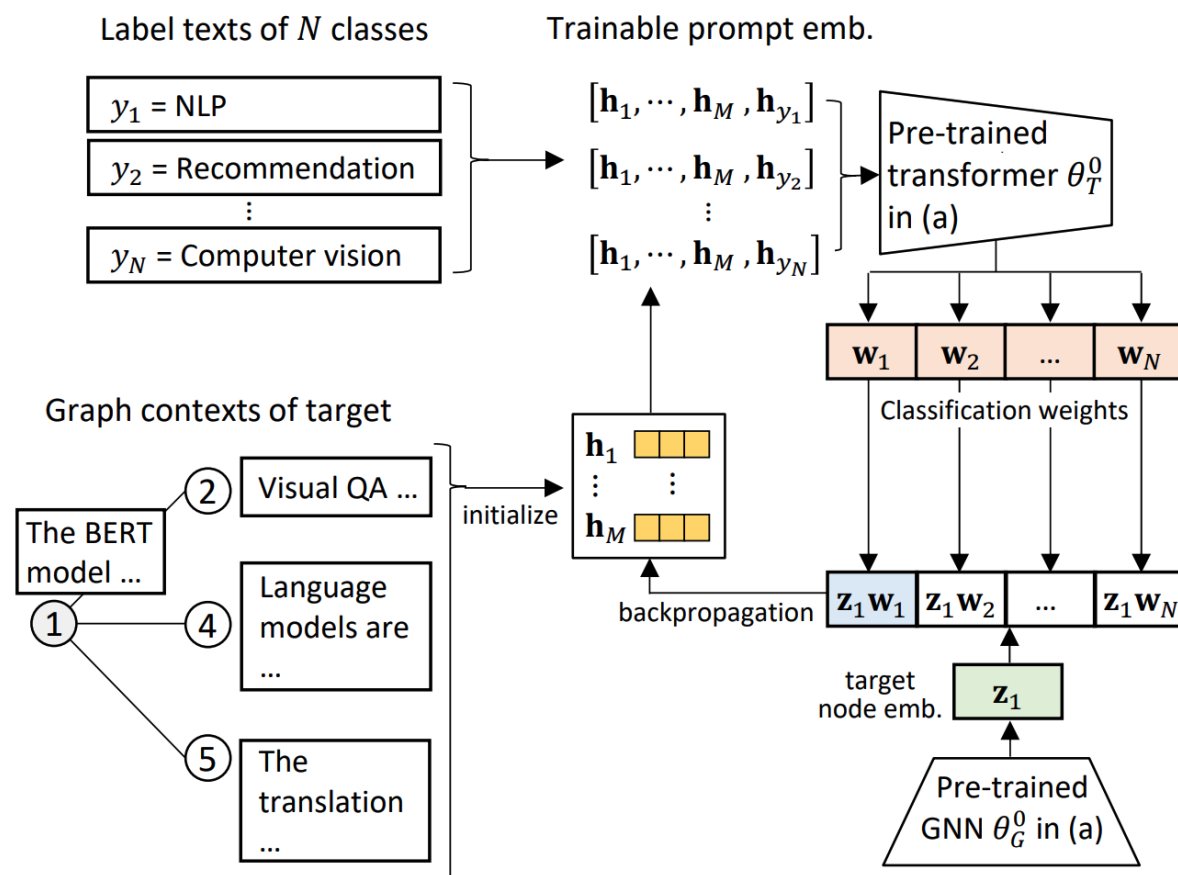
Graph-grounded pre-training and prompting (G2P2)

10

Zero-shot node classification with discrete prompts



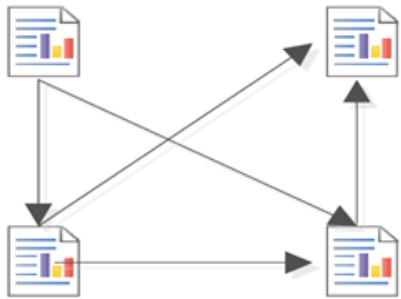
Few-shot node classification with continuous prompt tuning



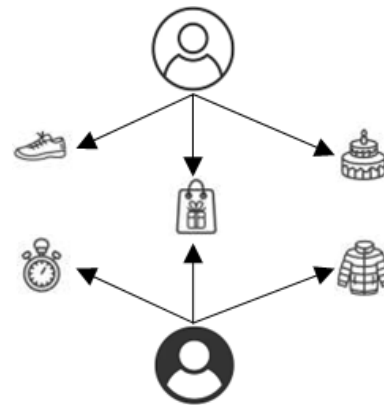
Datasets to evaluate G2P2

11

Dataset	Cora	Art	Industrial	M.I.
# Documents	25,120	1,615,902	1,260,053	905,453
# Links	182,280	4,898,218	3,101,670	2,692,734
# Avg. doc length	141.26	54.23	52.15	84.66
# Avg. node deg	7.26	3.03	2.46	2.97
# Classes	70	3,347	2,462	1,191



Cora is a collection of research papers with citation links



Art, Industrial and Music Instruments (M.I.) are three Amazon review datasets

Empirical performance of G2P2

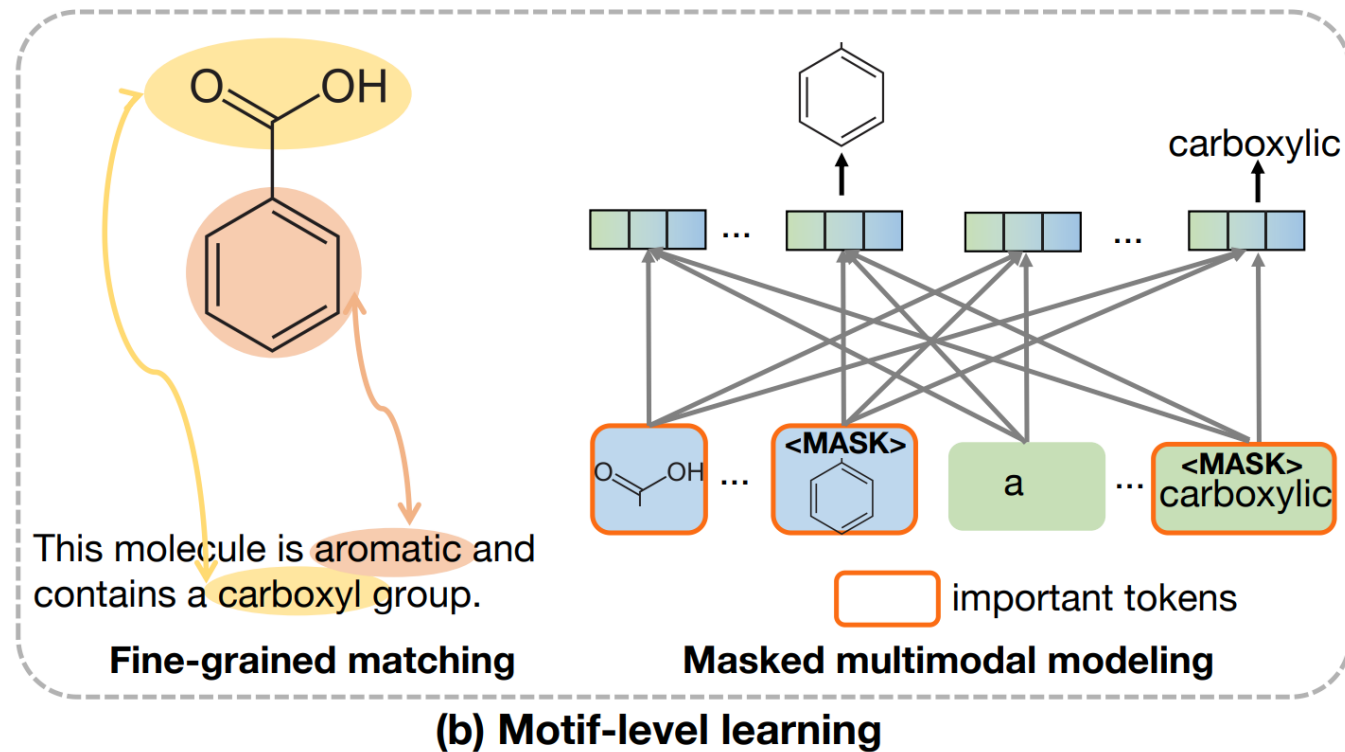
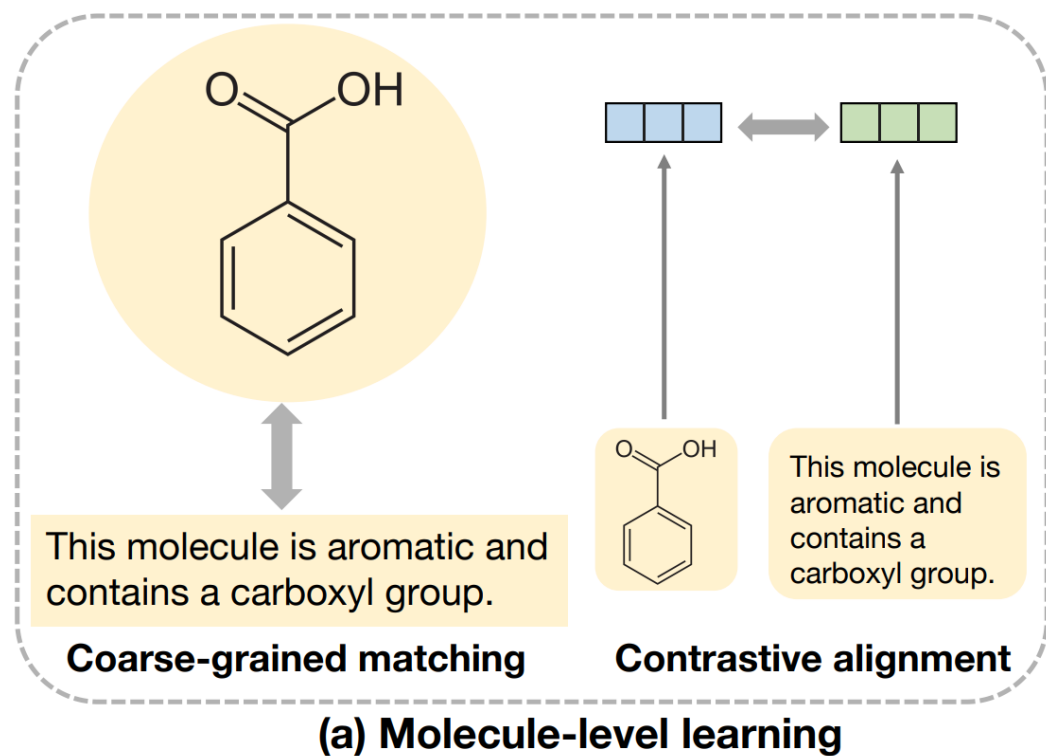
12

		Cora		Art		Industrial		M.I.	
		Accuracy	Macro-F1	Accuracy	Macro-F1	Accuracy	Macro-F1	Accuracy	Macro-F1
End-to-end GNN	GCN	41.15±2.41	34.50±2.23	22.47±1.78	15.45±1.14	21.08±0.45	15.23±0.29	22.54±0.82	16.26±0.72
	SAGE _{sup}	41.42±2.90	35.14±2.14	22.60±0.56	16.01±0.28	20.74±0.91	15.31±0.37	22.14±0.80	16.69±0.62
	TextGCN	59.78±1.88	55.85±1.50	43.47±1.02	32.20±1.30	53.60±0.70	45.97±0.49	46.26±0.91	38.75±0.78
Pre-trained GNN	GPT-GNN	76.72±2.02	72.23±1.17	65.15±1.37	52.79±0.83	62.13±0.65	54.47±0.67	67.97±2.49	59.89±2.51
	DGI	<u>78.42</u> ±1.39	<u>74.58</u> ±1.24	65.41±0.86	53.57±0.75	52.29±0.66	45.26±0.51	68.06±0.73	60.64±0.61
	SAGE _{self}	77.59±1.71	73.47±1.53	76.13±0.94	65.25±0.31	71.87±0.61	65.09±0.47	<u>77.70</u> ±0.48	<u>70.87</u> ±0.59
Pre-trained Transformers	BERT	37.86±5.31	32.78±5.01	46.39±1.05	37.07± 0.68	54.00±0.20	47.57±0.50	50.14±0.68	42.96±1.02
	BERT*	27.22±1.22	23.34±1.11	45.31±0.96	36.28±0.71	49.60±0.27	43.36±0.27	40.19±0.74	33.69±0.72
	RoBERTa	62.10±2.77	57.21±2.51	72.95±1.75	62.25±1.33	76.35±0.65	70.49±0.59	70.67±0.87	63.50±1.11
	RoBERTa*	67.42±4.35	62.72±3.02	74.47±1.00	63.35±1.09	77.08±1.02	71.44±0.87	74.61±1.08	67.78±0.95
Prompt tuning	P-Tuning v2	71.00±2.03	66.76±1.95	<u>76.86</u> ±0.59	<u>66.89</u> ±1.14	<u>79.65</u> ±0.38	<u>74.33</u> ±0.37	72.08±0.51	65.44±0.63
	G2P2-p	79.16±1.23	74.99±1.35	79.59±0.31	68.26±0.43	80.86±0.40	74.44±0.29	81.26±0.36	74.82±0.45
	G2P2	80.08* ±1.33	75.91* ±1.39	81.03* ±0.43	69.86* ±0.67	82.46* ±0.29	76.36* ±0.25	82.77* ±0.32	76.48* ±0.52
	(improv.)	(+2.12%)	(+1.78%)	(+5.43%)	(+4.44%)	(+3.53%)	(+2.7%)	(+6.53%)	(+7.92%)

G2P2 outperforms the best baseline (at that time) by around 3–7%.

Fine-grained graph-text integration

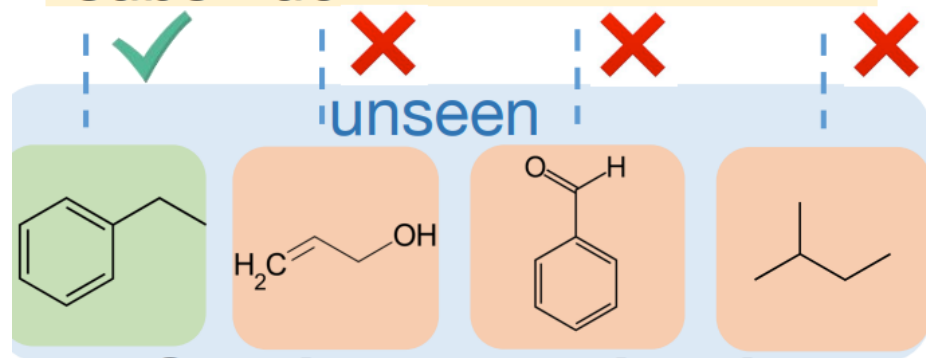
13



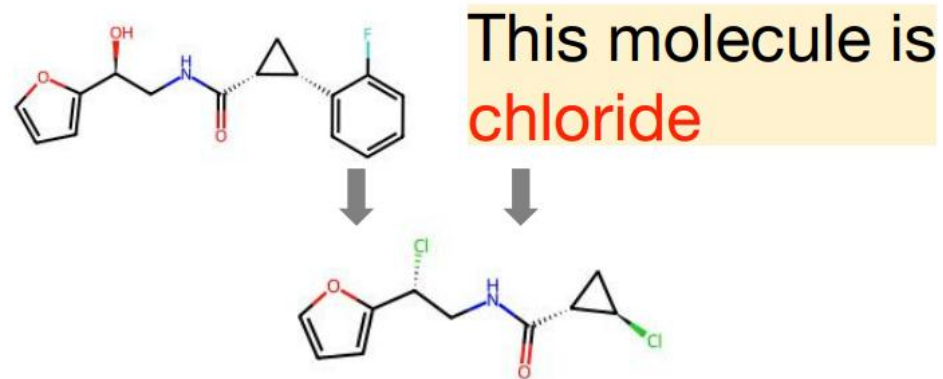
Why is fine-grained alignment important?

14

This molecule is an **alkyl-benzene** carrying an **ethyl** substituent.



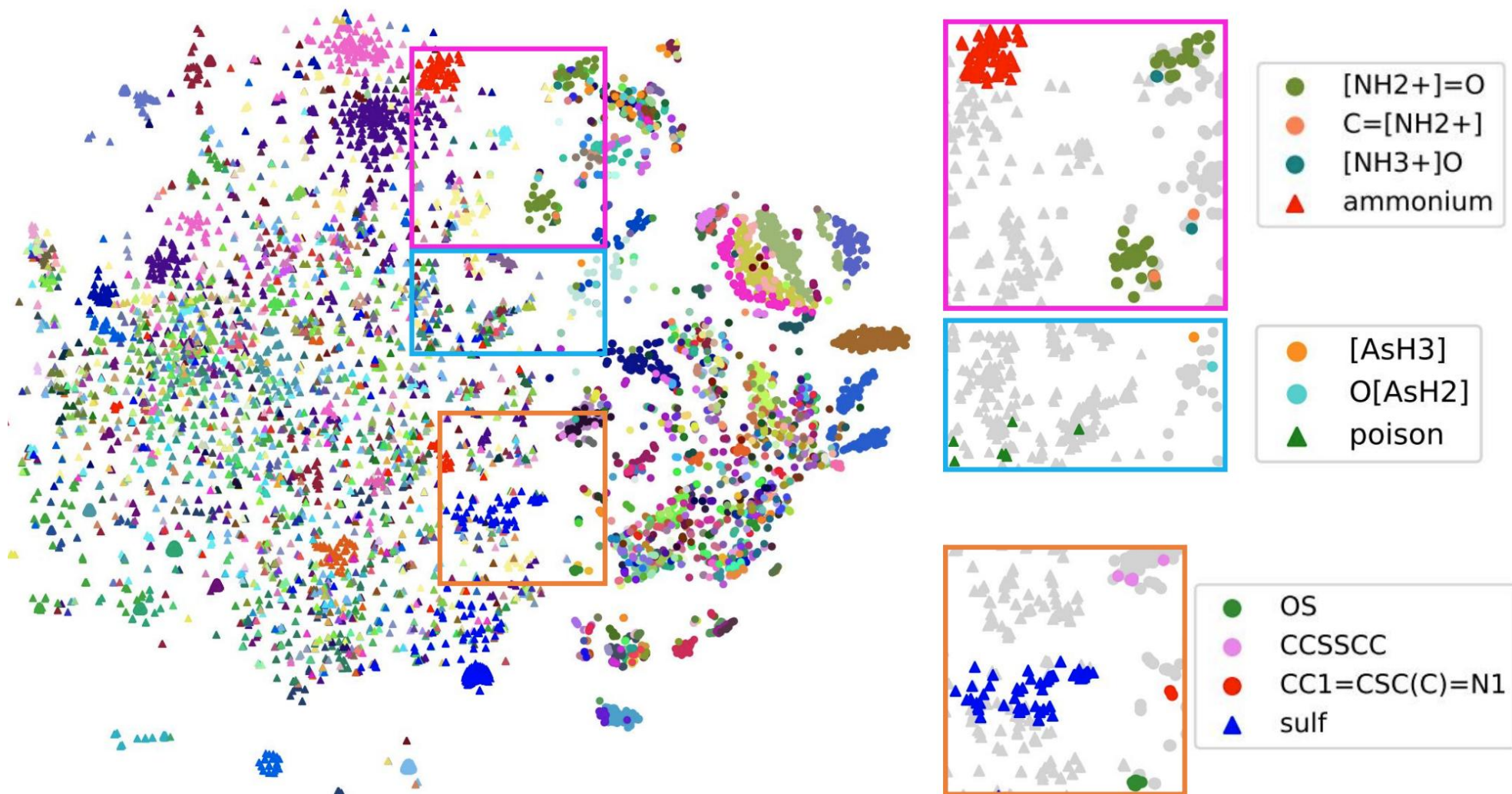
Graph-text retrieval



Text-based molecule editing

Visualization of learned word/motif embeddings

15



Outline

16

- Introduction: Graphs & text-attributed graphs
- Jointly training graph and textual data
- **Quantizing graphs into language tokens for LLMs**
- Conclusions

Integrating graph data in the era of LLMs

17

Graph verbalization

Instructor:

You are a brilliant graph master that can handle anything related to graphs like retrieval, detection and classification.

Graph description language:

```
<?xml version='1.0' encoding='utf-8'?>
<graphml xmlns="http://graphml.graphdrawing.org/xmlns">
  <key id="relation" for="edge" attr.name="relation" attr.type="string" />
  <key id="title" for="node" attr.name="title" attr.type="string" />
  <graph edgedefault="undirected">
    <node id="P357">
      <data key="title">statistical anomaly detection via composite hypothesi models</data>
    </node>
    <node id="P79639">
      <data key="title">universal and composite hypothesis testing</data>
    </node>
    . . . . .
    <edge source="P357" target="P79639">
      <data key="relation">reference</data>
    </edge>
    . . . . .
  </graph>
</graphml>
```

Context: XXXXXX

Query:

What is the clustering coefficient of node P357 ?

New Contexts:

Node P357 has 4 neighbors, where each of which are about anomaly detection with statsitital models. The whole graph contains 5 nodes and 10 edges and describes the citation relations.

Generate
New Contexts

LLMs

Generate
Final Output

Final Output:

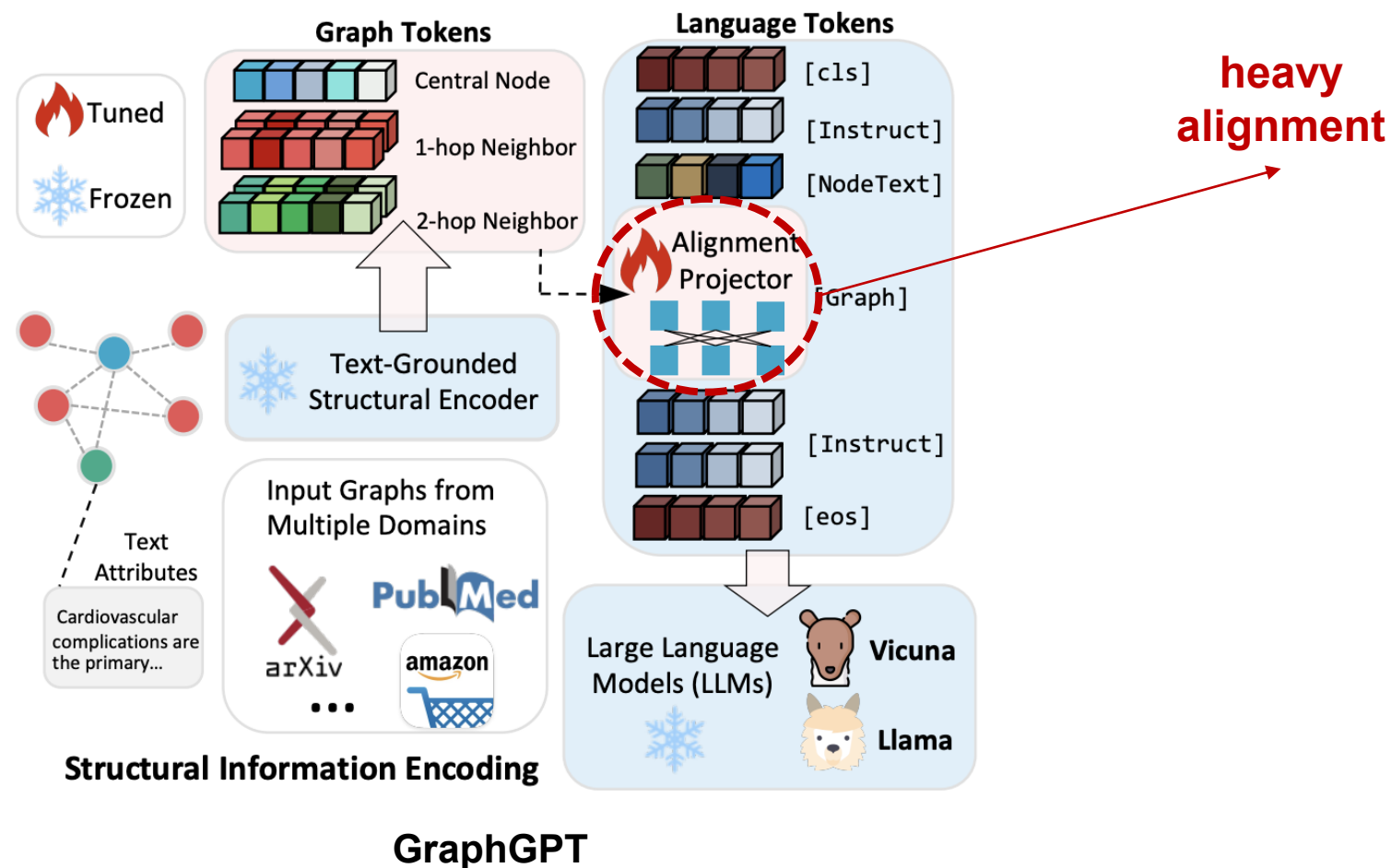
The clustering coefficient of a node is the ratio of the number of closed neighbors and all possible closed neighbors. To compute the clustering coefficient of node P357, we first know that node P357 has 4 neighbors

.....
that is there are totally 6 possible triangles between them. Among these neighbors, there are only 2 of them connected with each other, which forms one triangle. Thus the clustering coefficient of node P357 is $1 / 6 = 0.167$.

GPT4Graph

Integrating graph data in the era of LLMs

18



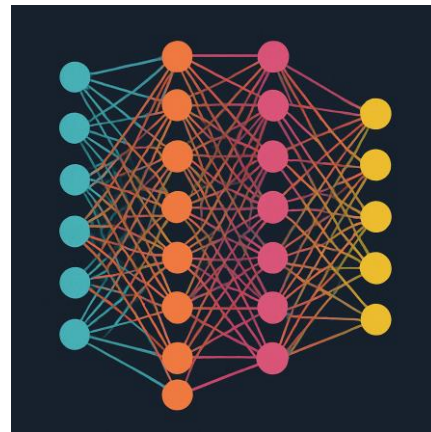
Structural-semantic gap

19

```
Instructor:
You are a brilliant graph master that can handle anything
related to graphs like retrieval, detection and classification.
Graph description language:
<?xml version='1.0' encoding='utf-8'?>
<graphml xmlns="http://graphml.graphdrawing.org/xmlns">
  <key id="relation" for="edge" attr.name="relation" attr.type="string" />
  <key id="title" for="node" attr.name="title" attr.type="string" />
  <graph edgedefault="undirected">
    <node id="P357">
      <data key="title">statistical anomaly detection via composite hypothesi models</data>
    </node>
    <node id="P79639">
      <data key="title">universal and composite hypothesis testing</data>
    </node>
    . . . . .
    <edge source="P357" target="P79639">
      <data key="relation">reference</data>
    </edge>
    . . . . .
  </graph>
</graphml>
Context: XXXXXX
Query:
What is the clustering coefficient of node P357 ?
```

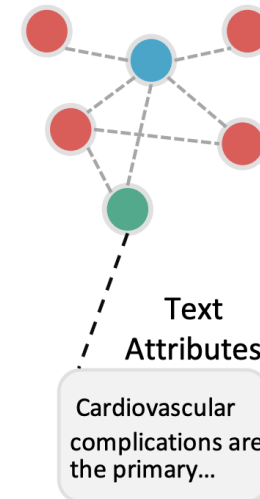
Graph Verbalization

Structural information loss



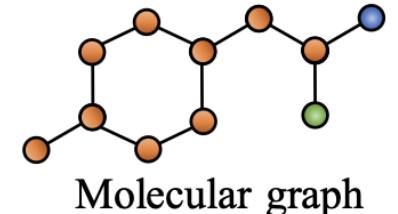
Projector-based Alignment

High computational cost



Transfer learning

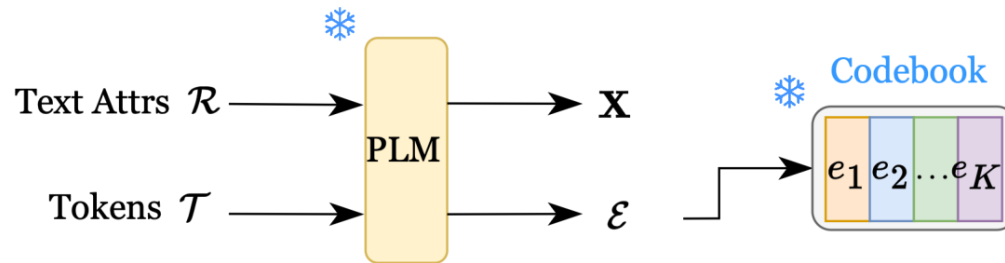
Poor generalization



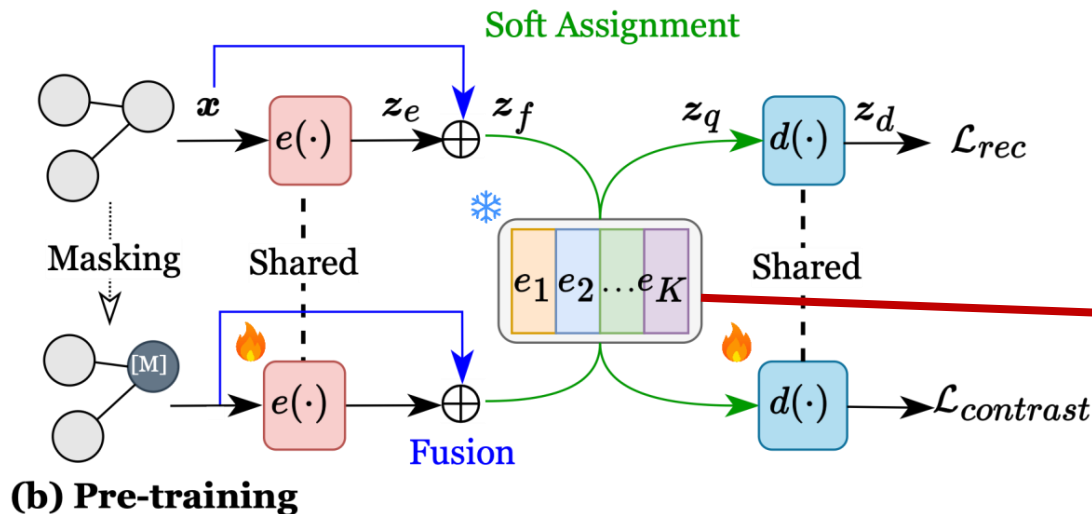
Continuous vs. Discrete
Graph embeddings \leftrightarrow LLM tokens

Soft Tokenization of Text-attributed Graphs (STAG)

20



(a) Codebook Construction



(b) Pre-training

AI	<div></div>	0.35
algorithm	<div></div>	0.25
complexity	<div></div>	0.20
computation	<div></div>	0.15
theory	<div></div>	0.05

Soft Tokenization of Text-attributed Graphs (STAG)

21

✓ With LLMs

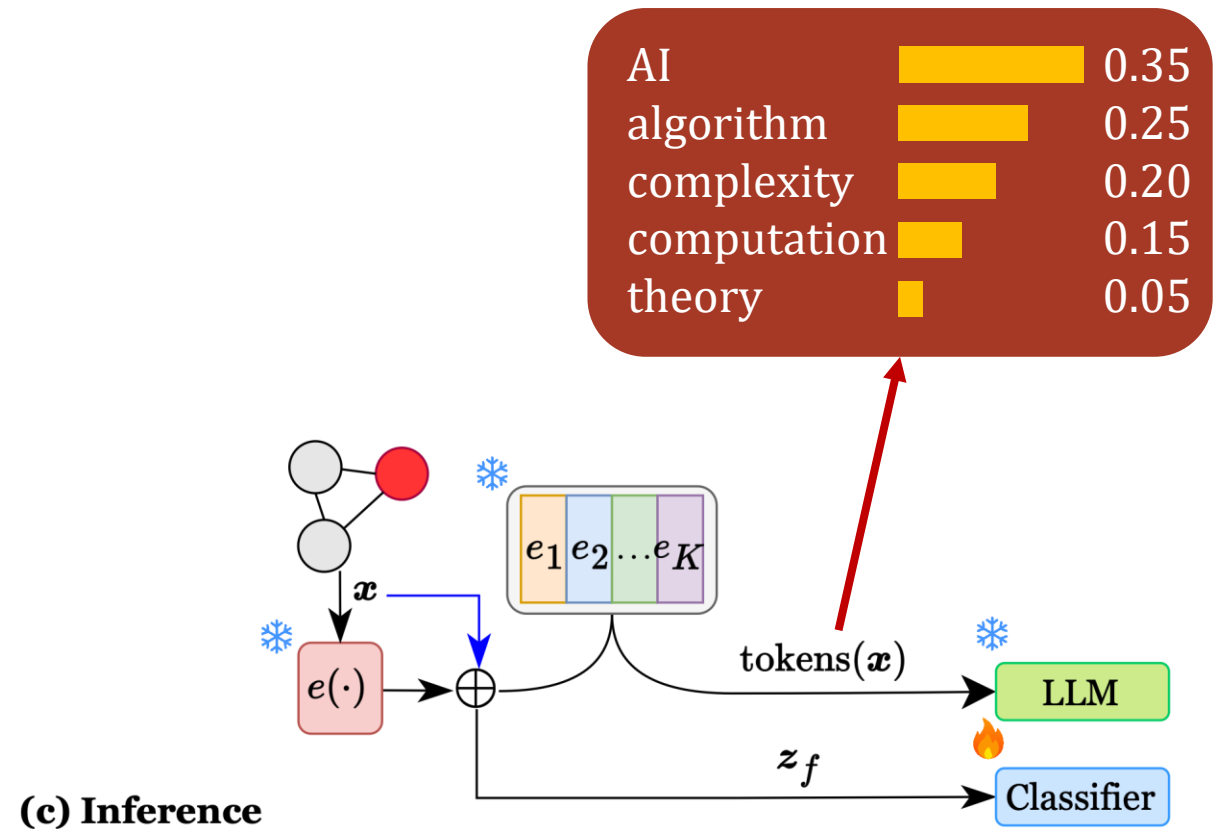
- Extract top-k tokens
- Few-shot: In-context learning
- Zero-shot: Direct LLM classification

✓ Without LLMs

- Linear probing on frozen embeddings

✓ Prompt Tuning

- Lightweight adaptation for domain transfer
- Supports both LLMs and without LLMs



Inference with LLMs

22

System Prompt: You are a node classifier. Given a list of tokens representing a node's features, predict its class from the following options: [Research Paper, Dataset, Software].

Few-shot examples: Node tokens: [research, methodology, experiment] Class: Research Paper

Node tokens: [benchmark, statistics, collection] Class: Dataset

Node tokens: [implementation, code, library] Class: Software

Test Node: Node tokens: [algorithm, computation, optimization] Predict the class:

Pre-train once, apply all

23

LLM	Cora Full	WikiCS	ogbn-arxiv	CiteSeer
LLaMA2-7B + PT	76.66 \pm 7.79 81.05 \pm 7.77	79.00 \pm 7.96 79.90 \pm 7.69	65.33 \pm 10.46 77.42 \pm 10.48	54.35 \pm 9.54 58.45 \pm 8.61
LLaMA2-13B + PT	77.62 \pm 8.67 81.95 \pm 7.06	79.80 \pm 7.30 80.45 \pm 7.66	69.38 \pm 8.83 77.75 \pm 9.01	54.60 \pm 8.79 57.30 \pm 9.20
Vicuna-7B + PT	74.12 \pm 6.47 80.77 \pm 6.75	80.30 \pm 7.02 80.10 \pm 7.39	64.84 \pm 9.38 76.95 \pm 9.43	49.25 \pm 6.72 52.25 \pm 8.23
Vicuna-13B + PT	77.76 \pm 8.58 81.38 \pm 7.65	79.35 \pm 7.98 79.25 \pm 7.50	66.03 \pm 9.34 75.65 \pm 9.59	52.25 \pm 6.39 53.00 \pm 8.16
LLaMA3-8B + PT	79.22 \pm 8.45 82.88 \pm 8.09	78.40 \pm 8.05 78.35 \pm 7.61	70.37 \pm 8.95 76.71 \pm 10.20	61.25 \pm 7.14 64.20 \pm 7.39
GPT-4o-mini + PT	79.25 \pm 8.42 83.04 \pm 7.84	81.05 \pm 6.80 81.90\pm6.16	71.32 \pm 9.13 77.51 \pm 9.58	61.90 \pm 7.22 65.90 \pm 7.04
GPT-4o + PT	81.40\pm7.41 83.28\pm7.06	81.45\pm7.10 81.60 \pm 7.19	72.75\pm8.83 78.85\pm9.74	62.95\pm6.61 65.90\pm7.03

- **Larger** models perform better
- **Newer** architectures show advantages
- **Prompt tuning** provides consistent gains

Outline

24

- Introduction: Graphs & text-attributed graphs
- Jointly training graph and textual data
- Quantizing graphs into language tokens for LLMs
- **Conclusions**

Conclusions

25

- Text-attributed graphs contain rich semantics
- Graph structures and semantics can be jointly pre-trained
- Quantizing graphs is promising for integration with LLMs

Acknowledgement

26



Zhihao Wen, Yuan Fang. Augmenting Low-Resource Text Classification with Graph-Grounded Pre-training and Prompting. *SIGIR* 2023.



Yibo Li, Yuan Fang, Mengmei Zhang, Chuan Shi. Advancing Molecular Graph-Text Pre-training via Fine-grained Alignment. *KDD* 2025.



Jianyuan Bo, Hao Wu, Yuan Fang. Quantizing Text-attributed Graphs for Semantic-Structural Integration. *KDD* 2025.



Jiawei Liu, Cheng Yang, Zhiyuan Lu, Junze Chen, Yibo Li, Mengmei Zhang, Ting Bai, Yuan Fang, Lichao Sun, Philip S. Yu, Chuan Shi. Graph Foundation Models: Concepts, Opportunities and Challenges. *TPAMI* 2025.

The research is made possible with support from the School of Computing and Information Systems, Singapore Management University. Full publications, codes and data are available at <http://www.yfang.site/>.