

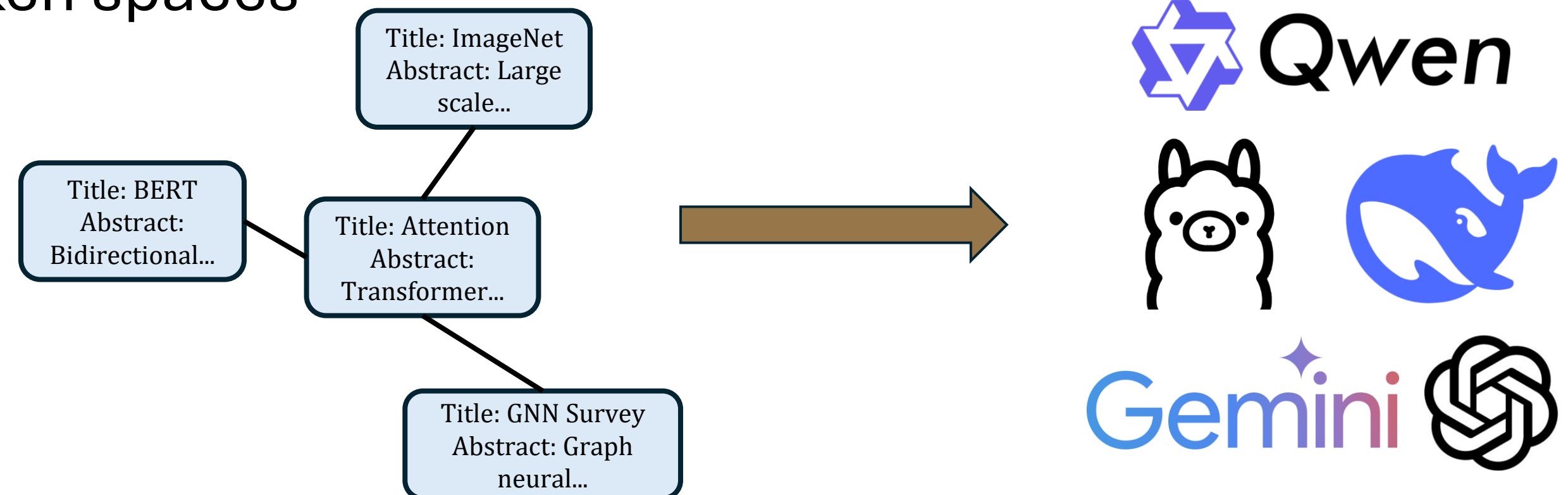
Quantizing Text-attributed Graphs for Semantic-Structural Integration

Introduction

Text-Attributed Graphs (TAGs)

+ Large Language Models

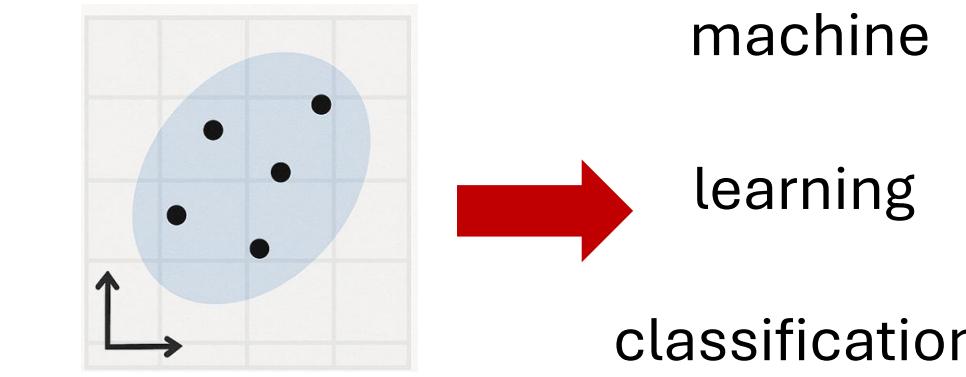
- TAGs contain rich semantic + structural information
- LLMs excel at semantic understanding but struggle with graph structures
- Need to bridge continuous graph embeddings \leftrightarrow discrete LLM token spaces



Motivation

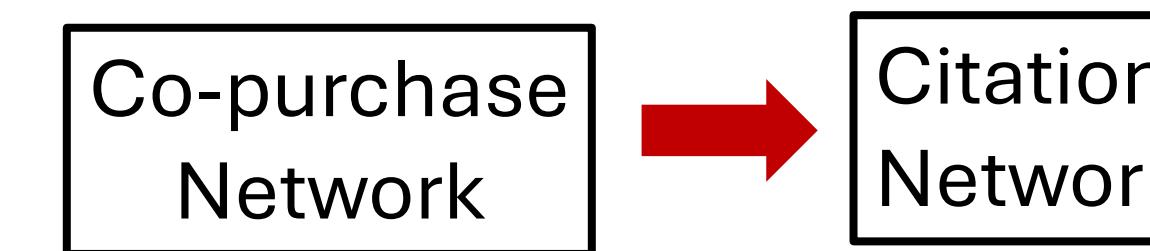
Problem 1: Semantic-Structural Integration

- Graph embeddings are continuous vectors
- LLMs work with discrete tokens



Problem 2: Transfer Learning Requirements

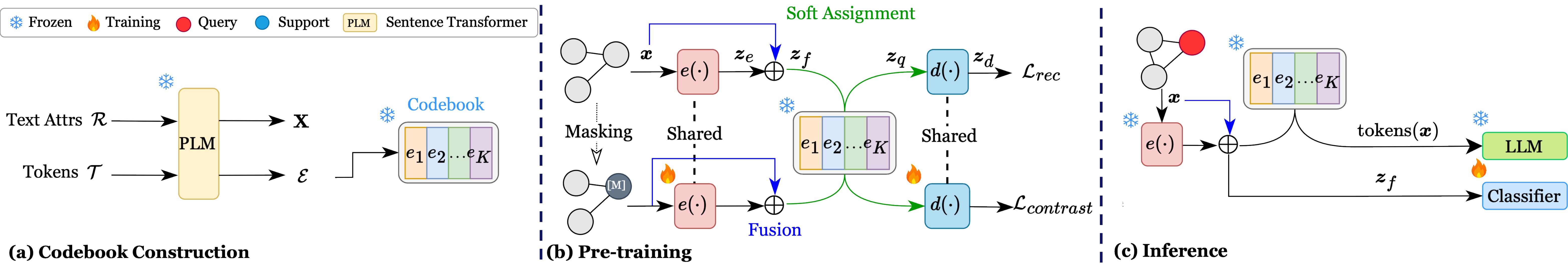
- Need labeled source data
- Limits adaptability across domains



Current Limitations

- X** Expensive alignment mechanisms between GNN and LLM
- X** Manual graph verbalization loses structural details
- X** Require labeled data for cross-dataset transfer

Proposed Method: STAG



Method

STAG: Soft Tokenization for Text-Attributed Graphs

Three Key Innovations

1 Semantic-Structural Fusion

- Fuse GNN embeddings + original text features

$$z_f = \phi \cdot \frac{\mathbf{W}_f z_e}{\|\mathbf{W}_f z_e\|_2} + \psi \cdot \frac{\mathbf{x}}{\|\mathbf{x}\|_2}$$

2 Soft Assignment Strategy

- Map nodes to token distributions (not single tokens)
- KL divergence guides quantization to semantic tokens
- Prevents overfitting, improves transferability



3 Dual-Branch Training

- Reconstruction branch: preserves semantics
- Contrastive branch: captures neighborhood structure
- No labeled data required!

Experiments

Few-shot Node Classification

Pre-train data	Method	LLM	Target data						
			Cora	Cora Full	CiteSeer	PubMed	WikiCS	ogbn-arxiv	
Same as target	GCN	\times	76.10 \pm 4.26	82.81 \pm 7.40	59.95 \pm 6.92	66.35 \pm 6.71	70.55 \pm 8.26	76.61 \pm 7.72	80.08 \pm 5.41
	GAT	\times	79.60 \pm 5.00	84.72 \pm 7.83	60.85 \pm 6.78	67.40 \pm 7.18	77.95 \pm 7.81	80.80 \pm 7.99	81.94 \pm 7.07
No pre-train	Raw Text	✓	63.40 \pm 9.07	71.66 \pm 7.84	62.10 \pm 5.35	85.00 \pm 5.48	77.15 \pm 6.92	54.74 \pm 9.21	87.58 \pm 5.48
	Raw Feat + Quantization	✓	54.85 \pm 6.28	73.74 \pm 8.10	56.40 \pm 5.48	48.30 \pm 7.78	73.40 \pm 8.19	63.75 \pm 9.98	65.84 \pm 9.96
	Raw Feat + Linear Probing	\times	70.25 \pm 7.22	81.29 \pm 7.47	63.00 \pm 6.72	68.30 \pm 6.28	78.05 \pm 7.47	83.05 \pm 7.40	77.53 \pm 7.16
Cora Full	DGI	\times	77.05 \pm 5.12	83.32 \pm 8.12	63.85 \pm 5.39	68.20 \pm 7.57	78.65 \pm 6.90	81.30 \pm 8.51	79.90 \pm 7.20
	GraphMAE2	\times	77.70 \pm 6.92	84.74 \pm 7.42	65.25 \pm 5.84	66.35 \pm 6.09	80.95 \pm 4.96	80.04 \pm 8.15	73.93 \pm 7.57
	GPPT	\times	27.16 \pm 7.61	67.90 \pm 12.72	28.66 \pm 7.60	21.53 \pm 10.91	29.00 \pm 8.08	36.92 \pm 10.32	24.32 \pm 5.13
	G2P2	\times	74.90 \pm 7.47	81.10 \pm 7.44	59.65 \pm 9.68	67.85 \pm 8.02	69.90 \pm 10.52	68.75 \pm 10.14	70.97 \pm 10.03
	Prodigy	\times	39.50 \pm 6.75	60.80 \pm 6.38	42.90 \pm 5.02	43.68 \pm 6.91	43.25 \pm 6.91	47.85 \pm 6.89	30.70 \pm 5.94
	OFA	\times	45.95 \pm 4.52	56.95 \pm 5.31	36.80 \pm 5.50	49.40 \pm 4.75	46.45 \pm 4.67	50.80 \pm 4.73	33.60 \pm 4.26
	STAG	✓	67.60 \pm 6.72	80.95 \pm 8.02	62.45 \pm 7.02	54.50 \pm 7.83	79.20 \pm 8.41	71.56 \pm 10.32	69.34 \pm 9.93
	+ Linear Probing	\times	78.50 \pm 5.62	86.04 \pm 6.70	66.70 \pm 5.36	69.00 \pm 6.31	84.05 \pm 5.72	82.99 \pm 8.10	79.62 \pm 7.12
	+ Prompt Tuning	✓	73.30 \pm 4.77	85.20 \pm 7.59	65.40 \pm 5.98	66.20 \pm 5.70	79.45 \pm 7.53	79.18 \pm 8.28	73.94 \pm 9.67
	+ Prompt Tuning*	\times	78.65 \pm 5.93	86.66 \pm 7.67	65.80 \pm 7.03	68.25 \pm 6.80	83.55 \pm 5.94	83.57 \pm 8.30	80.48 \pm 6.86

Zero-shot Node Classification

Pre-train data	Method	LLM	Target data			
			Cora	Cora Full	WikiCS	ogbn-arxiv
No pre-train	Raw Feat + Q	✓	47.10 \pm 5.98	60.33 \pm 10.88	70.40 \pm 8.88	25.48 \pm 5.54
	Raw Feat + C	\times	62.20 \pm 8.45	77.23 \pm 8.96	73.85 \pm 8.02	72.85 \pm 10.43
Cora Full	G2P2	\times	60.45 \pm 7.58	64.29 \pm 11.56	50.25 \pm 8.43	19.66 \pm 6.38
	OFA	\times	20.30 \pm 2.93	23.85 \pm 3.58	21.45 \pm 3.99	17.60 \pm 3.74
	STAG	✓	48.05 \pm 6.15	62.63 \pm 11.70	76.25 \pm 8.48	26.01 \pm 7.52
	STAG + C	\times	66.55 \pm 7.48	82.90 \pm 9.52	75.15 \pm 7.81	74.23 \pm 9.35

Inference with LLM

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:

Zero-shot Binary Link Prediction

Method	Cora	ogbn-products
LLaGA	87.35	92.99
STAG	63.00	92.65
STAG (non LLM)	93.20	96.85

N-way 5-shot Edge Classification

Method	WN18RR	FB1