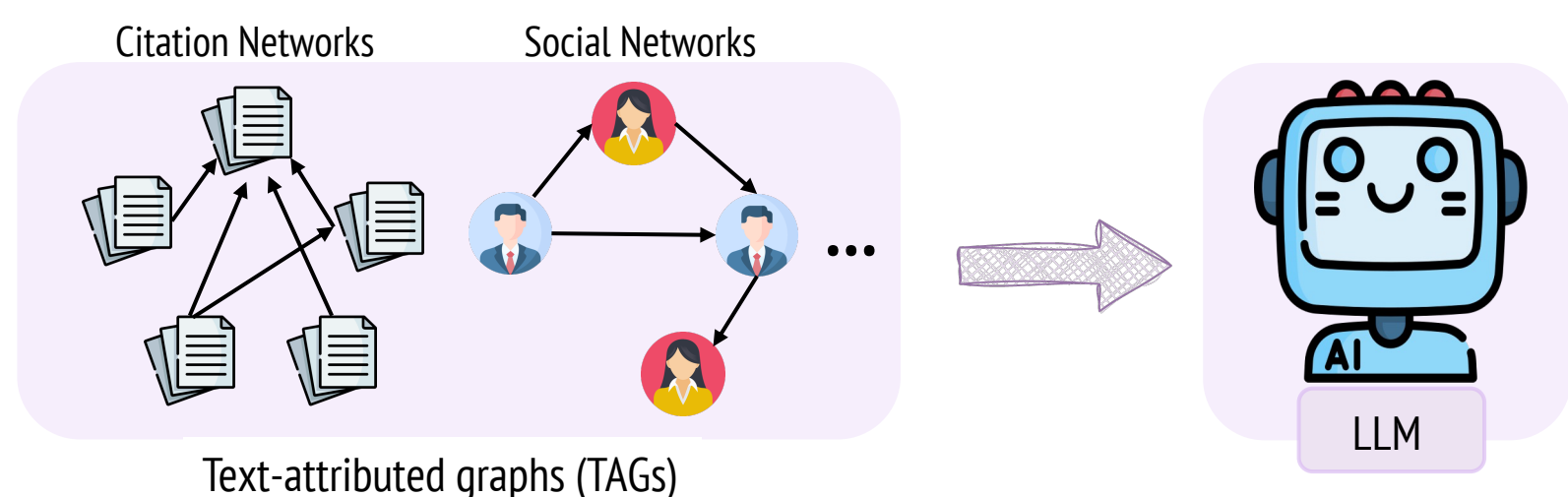


Background & Motivation

Introductoin

- LLMs are powerful for text-attributed graph (TAG) learning, leveraging superior semantic understanding to achieve strong node classification performance.
- LLM-as-Enhancers use LLMs to generate enriched node features or embeddings for downstream GNN learning.
- LLM-as-Predictors directly use LLMs to predict node classes as a text generation task, and are more flexible as they enable zero-shot transfer across datasets and label spaces.



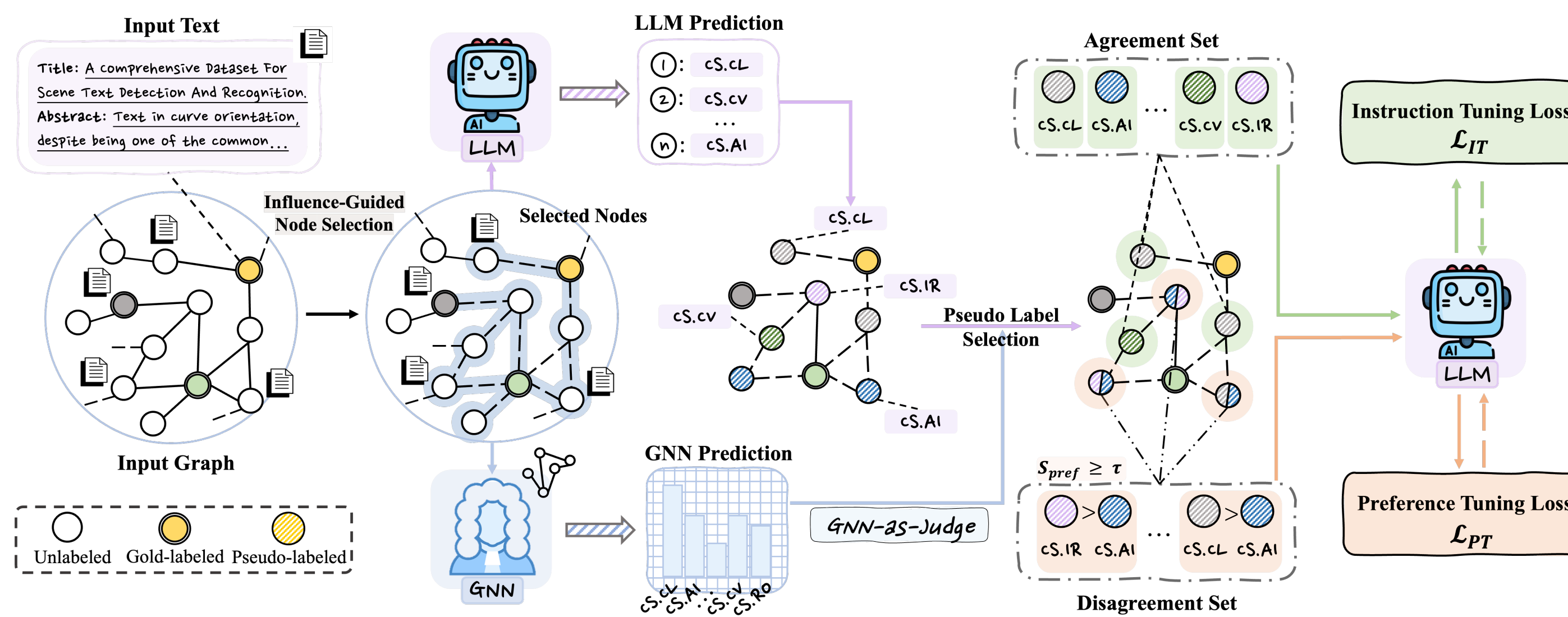
Problem

- LLM-as-Predictors require sufficient labeled nodes for fine-tuning – but real-world graphs are sparsely labeled.
- Under few-shot settings, LLMs lack structural inductive bias, leading to unreliable pseudo-labels and poor generalization.
- Challenge ①:** How to generate reliable pseudo-labels beyond LLM's own knowledge on graphs?
- Challenge ②:** How to mitigate label noise when fine-tuning LLMs with pseudo-labeled data?

Contributions

- We study LLM-as-Predictors for few-shot semi-supervised node classification on TAGs, a fundamental yet underexplored problem.
- We propose GNN-as-Judge, a novel framework that positions GNNs as judges to select reliable pseudo-labels by exploiting agreement and disagreement between GNNs and LLMs.
- We develop a weakly-supervised fine-tuning algorithm combining instruction tuning on agreement nodes and preference tuning on disagreement nodes, achieving state-of-the-art performance across multiple TAG benchmarks.

Method



Influence-Guided Node Selection for Pseudo Labeling

- Calculate estimated influence score for unlabeled nodes

$$\mathcal{I}\mathcal{S}(v_j) = \max_{v_i \in \mathcal{V}_{\text{train}}} \frac{|\mathcal{P}_{v_i, v_j}^*|}{(D_{GM}^*)^{h^*}}$$

Nodes closer to labeled nodes with shorter paths and lower-degree intermediaries receive stronger label signals

- Select Top-K most informative unlabeled nodes for pseudo labelling

$$\mathcal{V}_{\text{sel}} = \text{TopK}(\{\mathcal{I}\mathcal{S}(v_j)\}_{v_j \in \mathcal{V}_{\text{unlab}}}, K)$$

Agree/Disagree Node Set Selection with GNN Feedback

- Agreement Set $\mathcal{V}_{\text{agreed}}$:

$$\hat{y}^{\text{GNN}} = \hat{y}^{\text{LLM}}$$

Theorem 2: Intuitively, two models with different inductive biases are less likely to make the same mistake

- Disagreement Set $\mathcal{V}_{\text{disagreed}}$:

$$\hat{y}^{\text{GNN}} \neq \hat{y}^{\text{LLM}}$$

- Preference Score:

$$S_{\text{pref}}(v_i) = P_{\text{GNN}}(\hat{y}^{\text{GNN}} | v_i) - P_{\text{GNN}}(\hat{y}^{\text{LLM}} | v_i)$$

Select disagreed nodes with $S_{\text{pref}} \geq \tau$

Weakly-Supervised Fine-tuning for Graphs

- Agreement nodes use Instruction-Tuning loss:

$$\mathcal{L}_{IT}(\theta; x_i, y_i) = -\log p_{\theta}(y_i | x_i)$$

- Disagreement nodes use Preference-Tuning loss:

$$\mathcal{L}_{PT}(\theta; x_i, y_{w,i}, y_{l,i}) = -\log \sigma(g_{\theta}(x_i, y_{w,i}, y_{l,i}))$$

- Unified Objective:

$$\mathcal{L}(\theta) = \mathbb{E}[\mathcal{L}_{IT}] + \lambda \mathbb{E}[\mathcal{L}_{PT}]$$

IT consolidates correct predictions while PT gets learning signal from noisy labels with GNN feedback

Experiments

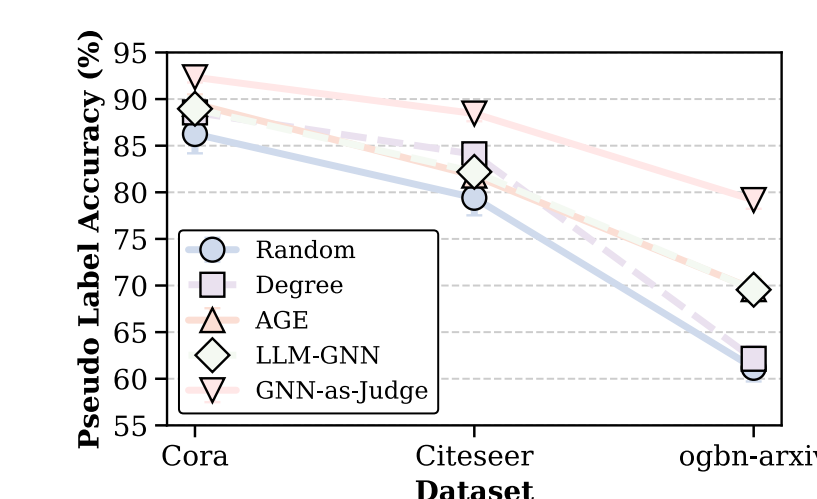
Main Results

Shot	Method	Cora	Citeseer	Pubmed	ogbn-arxiv	ogbn-products
3-shot	GCN	69.45±2.34	63.12±1.02	65.23±1.67	38.33±2.41	59.19±0.79
	SGC	67.21±1.25	63.07±0.95	65.34±1.18	39.16±1.49	57.42±1.12
	Zero-Shot	65.54±0.26	58.17±0.64	74.51±0.39	50.18±1.44	75.48±1.54
	Graph-CoT	63.02±0.77	47.23±1.06	86.22±2.47	49.67±1.23	74.15±1.83
	w. Neighbor	68.72±1.56	54.93±1.28	74.98±3.16	49.28±2.09	76.88±1.07
	GNN-as-Judge	77.89±1.28	73.59±0.64	87.12±0.89	62.21±1.45	81.02±1.23
5-shot	GCN	73.69±1.03	63.24±0.87	70.58±0.49	45.67±0.41	68.26±1.54
	SGC	72.48±0.35	62.08±0.77	70.74±0.94	49.89±1.23	66.78±1.56
	Zero-Shot	65.54±0.26	58.17±0.64	74.51±0.39	50.18±1.44	75.48±1.54
	Graph-CoT	63.02±0.77	47.23±1.06	86.22±2.47	49.67±1.23	74.15±1.83
	w. Neighbor	68.72±1.56	54.93±1.28	74.98±3.16	49.28±2.09	76.88±1.07
	GNN-as-Judge	79.54±0.39	74.39±1.63	87.49±1.23	66.76±0.83	81.93±2.21
10-shot	GCN	78.22±0.89	68.38±1.49	75.33±0.94	50.95±1.77	69.65±0.89
	SGC	78.49±0.37	67.44±0.60	74.98±1.91	51.89±1.23	67.91±0.48
	Zero-Shot	65.54±0.26	58.17±0.64	74.51±0.39	50.18±1.44	75.48±1.54
	Graph-CoT	63.02±0.77	47.23±1.06	86.22±2.47	49.67±1.23	74.15±1.83
	w. Neighbor	68.72±1.56	54.93±1.28	74.98±3.16	49.28±2.09	76.88±1.07
	GNN-as-Judge	80.71±0.83	74.62±1.35	90.17±1.69	67.88±1.03	82.48±1.56

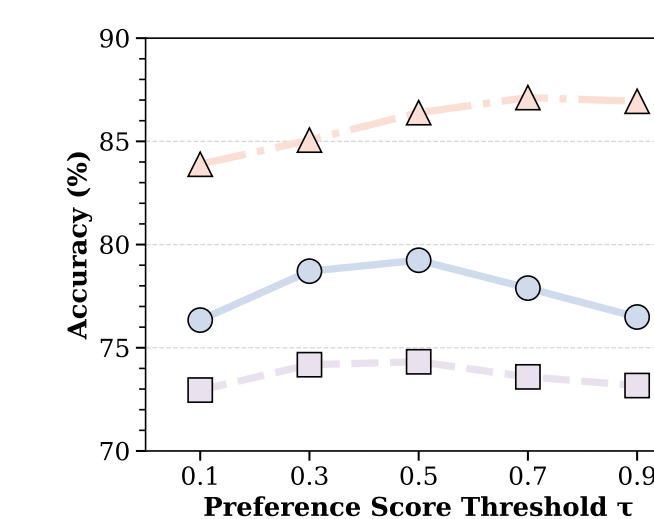
Zero-Shot Node Classification

Train → Test	Model	Accuracy
ogbn-arxiv	LLAGA	16.24±0.95
	GraphGPT	6.29±0.73
Cora	GNN-as-Judge	68.27±0.91
ogbn-arxiv	LLAGA	14.72±1.12
	GraphGPT	5.37±0.84
Citeseer	GNN-as-Judge	56.67±0.89
ogbn-arxiv	LLAGA	30.52±1.18
	GraphGPT	10.54±1.05
Pubmed	GNN-as-Judge	83.41±0.76

Pseudo-Label Selection



Sensitivity Analysis



Diversity Analysis

