

# When Do LLMs Help With Node Classification? A Comprehensive Analysis

Xixi Wu<sup>1</sup> Yifei Shen<sup>2</sup> Fangzhou Ge<sup>1</sup> Caihua Shan<sup>2</sup>  
Yizhu Jiao<sup>3</sup> Xiangguo Sun<sup>1</sup> Hong Cheng<sup>1</sup>

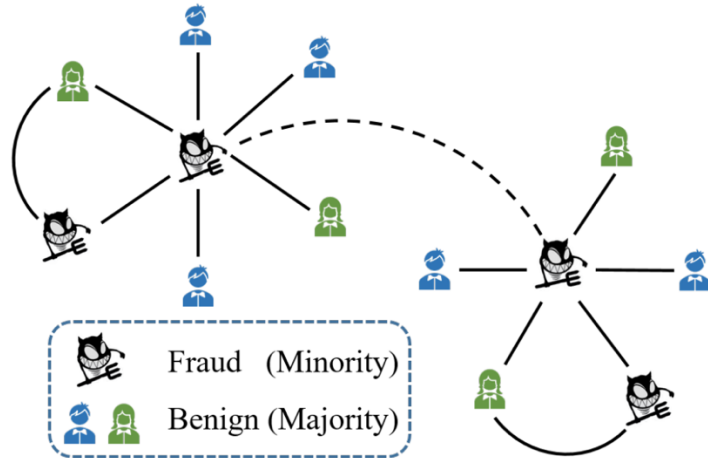
<sup>1</sup>The Chinese University of Hong Kong   <sup>2</sup>Microsoft Research Asia

<sup>3</sup>University of Illinois Urbana-Champaign

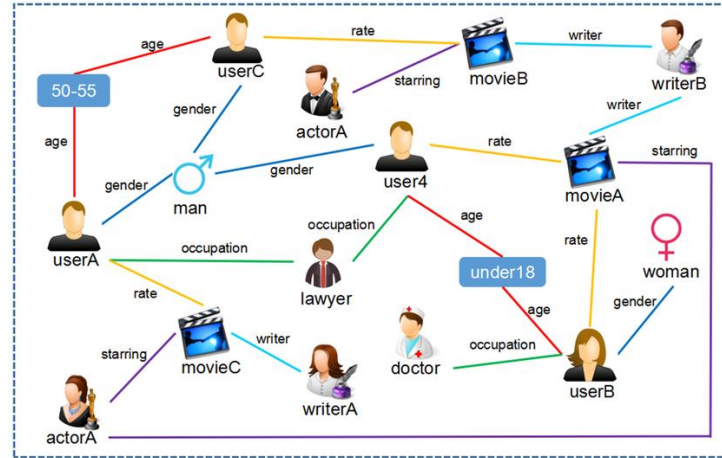


# Motivation

- Node classification is a fundamental task in graph analysis



*Fraud detection*



*User profiling*

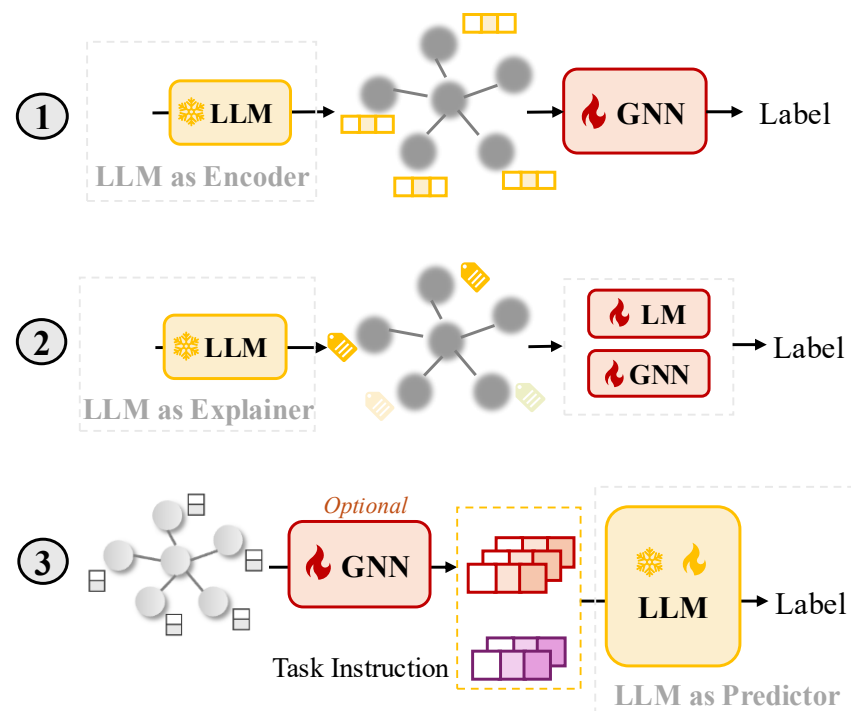
*Item tagging*

...

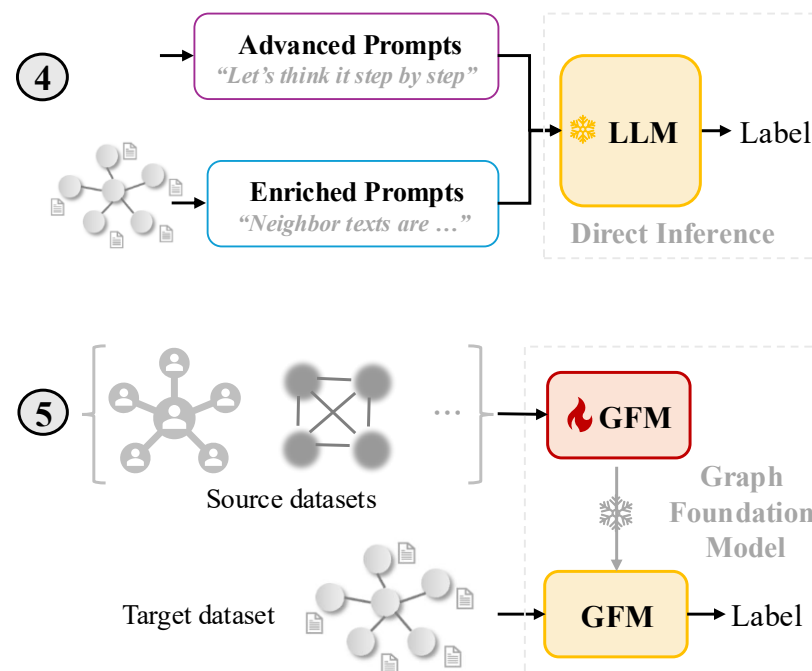
Even a marginal improvement in classification accuracy could result in **substantial financial profits**

# Motivation

- Leveraging LLMs for node classification has become popular



(a) Semi- / Supervised Settings



(b) Zero-shot Setting

*Superior semantic comprehension of LLMs overcome the limitations of shallow embeddings*

# Motivation

- Designing principles for LLM-based node classification algorithms remain elusive
  - For each algorithm category, what is the most suitable setting?
  - Under what scenarios, LLMs can surpass traditional LMs like BERT?
  - ...

## Unified Benchmark

Evaluate all methods using consistent dataloaders, learning paradigms, backbones, and implementation codebases

## Controlled Experiments

Consider comprehensive variables including learning paradigms, language model type & size, graph characteristics, etc

# Benchmark - LLMNodeBed

- Contains **14** datasets with varying scales, domain, and homophily
- Integrates **8** LLM-based algorithms, **8** classic methods
- Supports **3** different learning paradigms: semi-supervised, supervised, and zero-shot

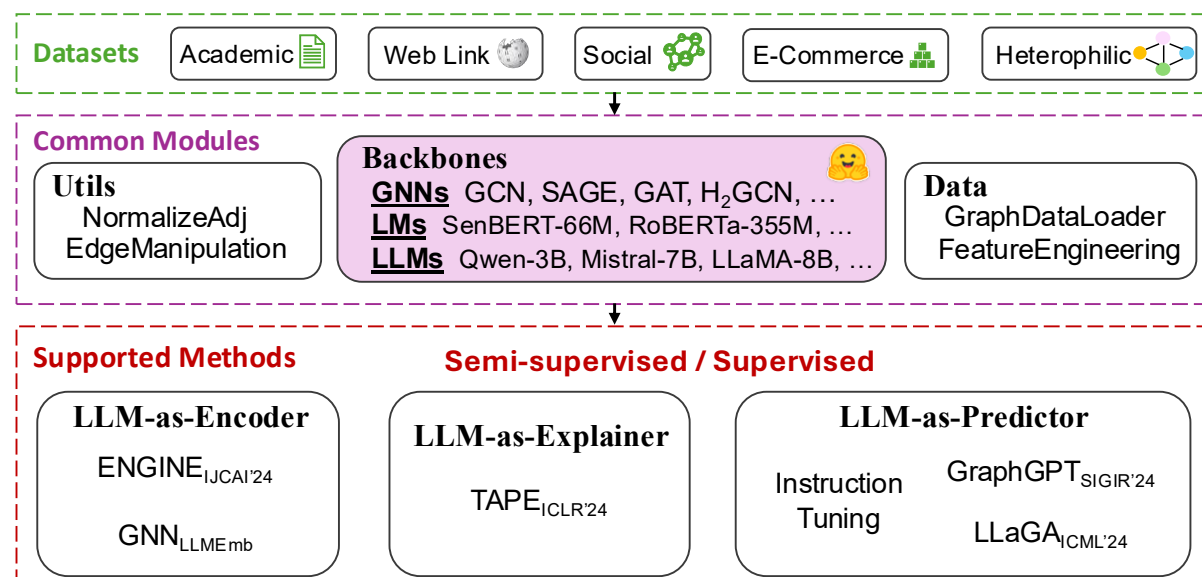


Table 1: Statistics of supported datasets in LLMNodeBed.

Statistics	Academic				Web Link	Social		E-Commerce			Heterophilic				Foundation Models
	Cora	Citeseer	Pubmed	arXiv	WikiCS	Instagram	Reddit	Books	Photo	Computer	Cornell	Texas	Wisconsin	Washington	
# Classes	7	6	3	40	10	2	2	12	12	10	5	5	5	5	Fine-tuned LLMs
# Nodes	2,708	3,186	19,717	169,343	11,701	11,339	33,434	41,551	48,362	87,229	191	187	265	229	
# Edges	5,429	4,277	44,338	1,166,243	215,863	144,010	198,448	358,574	500,928	721,081	292	310	510	394	Hallucination
Avg. # Token	183.4	210.0	446.5	239.8	629.9	56.2	197.3	337.0	201.5	123.1	594.6	453.2	639.1	469.0	
Homophily (%)	82.52	72.93	79.24	63.53	68.67	63.35	55.52	78.05	78.50	85.28	11.55	6.69	16.27	17.07	

# Experiments & Findings

- Semi-supervised & Supervised

Comparing Classic and LLM-based,  
*introducing LLMs to exploit textual  
information is useful*

Semi-supervised		Cora	Citeseer	Pubmed	WikiCS	Instagram	Reddit	Books	Photo	Computer	Avg.
Classic	GCN <sub>ShallowEmb</sub>	82.30 $\pm$ 0.19	70.55 $\pm$ 0.32	78.94 $\pm$ 0.27	79.86 $\pm$ 0.19	63.50 $\pm$ 0.11	61.44 $\pm$ 0.38	68.79 $\pm$ 0.46	69.25 $\pm$ 0.81	71.44 $\pm$ 1.19	71.79
	SAGE <sub>ShallowEmb</sub>	82.27 $\pm$ 0.37	69.56 $\pm$ 0.43	77.88 $\pm$ 0.44	79.67 $\pm$ 0.25	63.57 $\pm$ 0.10	56.65 $\pm$ 0.33	72.01 $\pm$ 0.33	78.50 $\pm$ 0.15	81.43 $\pm$ 0.27	73.50
	GAT <sub>ShallowEmb</sub>	81.30 $\pm$ 0.67	69.94 $\pm$ 0.74	78.49 $\pm$ 0.70	79.99 $\pm$ 0.65	63.56 $\pm$ 0.04	60.60 $\pm$ 1.17	74.35 $\pm$ 0.35	80.40 $\pm$ 0.45	83.39 $\pm$ 0.22	74.67
	SenBERT-66M	66.66 $\pm$ 1.42	60.52 $\pm$ 1.62	36.04 $\pm$ 2.92	77.77 $\pm$ 0.75	59.00 $\pm$ 1.17	56.05 $\pm$ 0.41	83.68 $\pm$ 0.19	73.89 $\pm$ 0.31	70.76 $\pm$ 0.15	64.93
	RoBERTa-355M	72.24 $\pm$ 1.14	66.68 $\pm$ 2.03	42.32 $\pm$ 1.56	76.81 $\pm$ 1.04	63.52 $\pm$ 0.44	59.27 $\pm$ 0.34	84.62 $\pm$ 0.16	74.79 $\pm$ 1.13	72.31 $\pm$ 0.37	68.06
	GLEM	81.30 $\pm$ 0.88	68.80 $\pm$ 2.46	81.70 $\pm$ 1.07	76.43 $\pm$ 0.55	60.25 $\pm$ 3.66	55.13 $\pm$ 1.41	83.28 $\pm$ 0.39	76.93 $\pm$ 0.49	80.46 $\pm$ 1.45	73.81
Encoder	GCN <sub>LLMEmb</sub>	83.33 $\pm$ 0.75	71.39 $\pm$ 0.90	78.71 $\pm$ 0.45	80.94 $\pm$ 0.16	67.49 $\pm$ 0.43	68.65 $\pm$ 0.75	83.03 $\pm$ 0.34	84.84 $\pm$ 0.47	88.22 $\pm$ 0.16	78.51
	ENGINE	84.22 $\pm$ 0.46	72.14 $\pm$ 0.74	77.84 $\pm$ 0.27	80.94 $\pm$ 0.19	67.14 $\pm$ 0.46	69.67 $\pm$ 0.16	82.89 $\pm$ 0.14	84.33 $\pm$ 0.57	86.42 $\pm$ 0.23	78.40
Explainer	TAPE	84.04 $\pm$ 0.24	71.87 $\pm$ 0.35	78.61 $\pm$ 1.23	81.94 $\pm$ 0.16	66.07 $\pm$ 0.10	62.43 $\pm$ 0.47	84.92 $\pm$ 0.26	86.46 $\pm$ 0.12	89.52 $\pm$ 0.04	78.43
Predictor	LLM <sub>IT</sub>	67.00 $\pm$ 0.16	54.26 $\pm$ 0.22	80.99 $\pm$ 0.43	75.02 $\pm$ 0.16	41.83 $\pm$ 0.47	54.09 $\pm$ 1.02	80.92 $\pm$ 1.38	71.28 $\pm$ 1.81	66.99 $\pm$ 2.02	65.76
	GraphGPT	64.72 $\pm$ 1.50	64.58 $\pm$ 1.55	70.34 $\pm$ 2.27	70.71 $\pm$ 0.37	62.88 $\pm$ 2.14	58.25 $\pm$ 0.37	81.13 $\pm$ 1.52	77.48 $\pm$ 0.78	80.10 $\pm$ 0.76	70.02
	LLaGA	78.94 $\pm$ 1.14	62.61 $\pm$ 3.63	65.91 $\pm$ 2.09	76.47 $\pm$ 2.20	65.84 $\pm$ 0.72	70.10 $\pm$ 0.38	83.47 $\pm$ 0.45	84.44 $\pm$ 0.90	87.82 $\pm$ 0.53	75.07

Supervised		Cora	Citeseer	Pubmed	arXiv	WikiCS	Instagram	Reddit	Books	Photo	Computer	Avg.
Classic	GCN <sub>ShallowEmb</sub>	87.41 $\pm$ 2.08	75.74 $\pm$ 1.20	89.01 $\pm$ 0.59	71.39 $\pm$ 0.28	83.67 $\pm$ 0.63	63.94 $\pm$ 0.61	65.07 $\pm$ 0.38	76.94 $\pm$ 0.26	73.34 $\pm$ 1.34	77.16 $\pm$ 3.80	76.37
	SAGE <sub>ShallowEmb</sub>	87.44 $\pm$ 1.74	74.96 $\pm$ 1.20	90.47 $\pm$ 0.25	71.21 $\pm$ 0.18	84.86 $\pm$ 0.91	64.14 $\pm$ 0.47	61.52 $\pm$ 0.60	79.40 $\pm$ 0.45	84.59 $\pm$ 0.32	87.77 $\pm$ 0.34	78.64
	GAT <sub>ShallowEmb</sub>	86.68 $\pm$ 1.12	73.73 $\pm$ 0.94	88.25 $\pm$ 0.47	71.57 $\pm$ 0.25	83.94 $\pm$ 0.61	64.93 $\pm$ 0.75	64.16 $\pm$ 1.05	80.61 $\pm$ 0.49	84.84 $\pm$ 0.69	88.32 $\pm$ 0.24	78.70
	SenBERT-66M	79.61 $\pm$ 1.40	74.06 $\pm$ 1.26	94.47 $\pm$ 0.33	72.66 $\pm$ 0.24	86.51 $\pm$ 0.86	60.11 $\pm$ 0.93	58.70 $\pm$ 0.54	85.99 $\pm$ 0.58	77.72 $\pm$ 0.35	74.22 $\pm$ 0.21	76.40
	RoBERTa-355M	83.17 $\pm$ 0.84	75.90 $\pm$ 1.69	94.84 $\pm$ 0.06	74.12 $\pm$ 0.12	87.47 $\pm$ 0.83	63.75 $\pm$ 1.13	60.61 $\pm$ 0.24	86.65 $\pm$ 0.38	79.45 $\pm$ 0.37	75.76 $\pm$ 0.30	78.17
	GLEM	86.81 $\pm$ 1.19	73.24 $\pm$ 1.55	93.98 $\pm$ 0.32	73.55 $\pm$ 0.22	79.81 $\pm$ 0.45	67.39 $\pm$ 1.73	53.11 $\pm$ 2.96	83.98 $\pm$ 0.97	78.16 $\pm$ 0.45	81.63 $\pm$ 0.46	77.17
Encoder	GCN <sub>LLMEmb</sub>	88.15 $\pm$ 1.79	76.45 $\pm$ 1.19	88.38 $\pm$ 0.68	74.39 $\pm$ 0.31	84.78 $\pm$ 0.86	68.27 $\pm$ 0.45	70.65 $\pm$ 0.75	84.23 $\pm$ 0.20	86.07 $\pm$ 0.20	89.52 $\pm$ 0.31	81.09
	ENGINE	87.00 $\pm$ 1.60	75.82 $\pm$ 1.52	90.08 $\pm$ 0.16	74.69 $\pm$ 0.36	85.44 $\pm$ 0.53	68.87 $\pm$ 0.25	71.21 $\pm$ 0.77	84.09 $\pm$ 0.09	86.98 $\pm$ 0.06	89.05 $\pm$ 0.13	81.32
Explainer	TAPE	88.05 $\pm$ 1.76	76.45 $\pm$ 1.60	93.00 $\pm$ 0.13	74.96 $\pm$ 0.14	87.11 $\pm$ 0.66	68.11 $\pm$ 0.54	66.22 $\pm$ 0.83	85.95 $\pm$ 0.59	87.72 $\pm$ 0.28	90.46 $\pm$ 0.18	81.80
Predictor	LLM <sub>IT</sub>	71.93 $\pm$ 1.47	60.97 $\pm$ 3.97	94.16 $\pm$ 0.19	76.08	80.61 $\pm$ 0.47	44.20 $\pm$ 3.06	58.30 $\pm$ 0.48	84.80 $\pm$ 0.13	78.27 $\pm$ 0.54	74.51 $\pm$ 0.53	72.38
	GraphGPT	82.29 $\pm$ 0.26	74.67 $\pm$ 1.15	93.54 $\pm$ 0.22	75.15 $\pm$ 0.14	82.54 $\pm$ 0.23	67.00 $\pm$ 1.22	60.72 $\pm$ 1.47	85.38 $\pm$ 0.72	84.46 $\pm$ 0.36	86.78 $\pm$ 1.14	79.25
	LLaGA	87.55 $\pm$ 1.15	76.73 $\pm$ 1.70	90.28 $\pm$ 0.91	74.49 $\pm$ 0.23	84.03 $\pm$ 1.10	69.16 $\pm$ 0.72	71.06 $\pm$ 0.38	85.56 $\pm$ 0.30	87.62 $\pm$ 0.30	90.41 $\pm$ 0.12	81.69

# Experiments & Findings

- Semi-supervised & Supervised

Comparing Semi-supervised and Supervised,  
**LLMs can bring greater improvements in  
semi-supervised settings than supervised**

Semi-supervised		Cora	Citeseer	Pubmed	WikiCS	Instagram	Reddit	Books	Photo	Computer	Avg.
Classic	GCN <sub>ShallowEmb</sub>	82.30 $\pm$ 0.19	70.55 $\pm$ 0.32	78.94 $\pm$ 0.27	79.86 $\pm$ 0.19	63.50 $\pm$ 0.11	61.44 $\pm$ 0.38	68.79 $\pm$ 0.46	69.25 $\pm$ 0.81	71.44 $\pm$ 1.19	71.79
	SAGE <sub>ShallowEmb</sub>	82.27 $\pm$ 0.37	69.56 $\pm$ 0.43	77.88 $\pm$ 0.44	79.67 $\pm$ 0.25	63.57 $\pm$ 0.10	56.65 $\pm$ 0.33	72.01 $\pm$ 0.33	78.50 $\pm$ 0.15	81.43 $\pm$ 0.27	73.50
	GAT <sub>ShallowEmb</sub>	81.30 $\pm$ 0.67	69.94 $\pm$ 0.74	78.49 $\pm$ 0.70	79.99 $\pm$ 0.65	63.56 $\pm$ 0.04	60.60 $\pm$ 1.17	74.35 $\pm$ 0.35	80.40 $\pm$ 0.45	83.39 $\pm$ 0.22	74.67
	SenBERT-66M	66.66 $\pm$ 1.42	60.52 $\pm$ 1.62	36.04 $\pm$ 2.92	77.77 $\pm$ 0.75	59.00 $\pm$ 1.17	56.05 $\pm$ 0.41	83.68 $\pm$ 0.19	73.89 $\pm$ 0.31	70.76 $\pm$ 0.15	64.93
	RoBERTa-355M	72.24 $\pm$ 1.14	66.68 $\pm$ 2.03	42.32 $\pm$ 1.56	76.81 $\pm$ 1.04	63.52 $\pm$ 0.44	59.27 $\pm$ 0.34	84.62 $\pm$ 0.16	74.79 $\pm$ 1.13	72.31 $\pm$ 0.37	68.06
	GLEM	81.30 $\pm$ 0.88	68.80 $\pm$ 2.46	81.70 $\pm$ 1.07	76.43 $\pm$ 0.55	60.25 $\pm$ 3.66	55.13 $\pm$ 1.41	83.28 $\pm$ 0.39	76.93 $\pm$ 0.49	80.46 $\pm$ 1.45	73.81
Encoder	GCN <sub>LLMEmb</sub>	83.33 $\pm$ 0.75	71.39 $\pm$ 0.90	78.71 $\pm$ 0.45	80.94 $\pm$ 0.16	67.49 $\pm$ 0.43	68.65 $\pm$ 0.75	83.03 $\pm$ 0.34	84.84 $\pm$ 0.47	88.22 $\pm$ 0.16	78.51
	ENGINE	84.22 $\pm$ 0.46	72.14 $\pm$ 0.74	77.84 $\pm$ 0.27	80.94 $\pm$ 0.19	67.14 $\pm$ 0.46	69.67 $\pm$ 0.16	82.89 $\pm$ 0.14	84.33 $\pm$ 0.57	86.42 $\pm$ 0.23	78.40
Explainer	TAPE	84.04 $\pm$ 0.24	71.87 $\pm$ 0.35	78.61 $\pm$ 1.23	81.94 $\pm$ 0.16	66.07 $\pm$ 0.10	62.43 $\pm$ 0.47	84.92 $\pm$ 0.26	86.46 $\pm$ 0.12	89.52 $\pm$ 0.04	78.43
Predictor	LLM <sub>IT</sub>	67.00 $\pm$ 0.16	54.26 $\pm$ 0.22	80.99 $\pm$ 0.43	75.02 $\pm$ 0.16	41.83 $\pm$ 0.47	54.09 $\pm$ 1.02	80.92 $\pm$ 1.38	71.28 $\pm$ 1.81	66.99 $\pm$ 2.02	65.76
	GraphGPT	64.72 $\pm$ 1.50	64.58 $\pm$ 1.55	70.34 $\pm$ 2.27	70.71 $\pm$ 0.37	62.88 $\pm$ 2.14	58.25 $\pm$ 0.37	81.13 $\pm$ 1.52	77.48 $\pm$ 0.78	80.10 $\pm$ 0.76	70.02
	LLaGA	78.94 $\pm$ 1.14	62.61 $\pm$ 3.63	65.91 $\pm$ 2.09	76.47 $\pm$ 2.20	65.84 $\pm$ 0.72	70.10 $\pm$ 0.38	83.47 $\pm$ 0.45	84.44 $\pm$ 0.90	87.82 $\pm$ 0.53	75.07

Supervised		Cora	Citeseer	Pubmed	arXiv	WikiCS	Instagram	Reddit	Books	Photo	Computer	Avg.
Classic	GCN <sub>ShallowEmb</sub>	87.41 $\pm$ 2.08	75.74 $\pm$ 1.20	89.01 $\pm$ 0.59	71.39 $\pm$ 0.28	83.67 $\pm$ 0.63	63.94 $\pm$ 0.61	65.07 $\pm$ 0.38	76.94 $\pm$ 0.26	73.34 $\pm$ 1.34	77.16 $\pm$ 3.80	76.37
	SAGE <sub>ShallowEmb</sub>	87.44 $\pm$ 1.74	74.96 $\pm$ 1.20	90.47 $\pm$ 0.25	71.21 $\pm$ 0.18	84.86 $\pm$ 0.91	64.14 $\pm$ 0.47	61.52 $\pm$ 0.60	79.40 $\pm$ 0.45	84.59 $\pm$ 0.32	87.77 $\pm$ 0.34	78.64
	GAT <sub>ShallowEmb</sub>	86.68 $\pm$ 1.12	73.73 $\pm$ 0.94	88.25 $\pm$ 0.47	71.57 $\pm$ 0.25	83.94 $\pm$ 0.61	64.93 $\pm$ 0.75	64.16 $\pm$ 1.05	80.61 $\pm$ 0.49	84.84 $\pm$ 0.69	88.32 $\pm$ 0.24	78.70
	SenBERT-66M	79.61 $\pm$ 1.40	74.06 $\pm$ 1.26	94.47 $\pm$ 0.33	72.66 $\pm$ 0.24	86.51 $\pm$ 0.86	60.11 $\pm$ 0.93	58.70 $\pm$ 0.54	85.99 $\pm$ 0.58	77.72 $\pm$ 0.35	74.22 $\pm$ 0.21	76.40
	RoBERTa-355M	83.17 $\pm$ 0.84	75.90 $\pm$ 1.69	94.84 $\pm$ 0.06	74.12 $\pm$ 0.12	87.47 $\pm$ 0.83	63.75 $\pm$ 1.13	60.61 $\pm$ 0.24	86.65 $\pm$ 0.38	79.45 $\pm$ 0.37	75.76 $\pm$ 0.30	78.17
	GLEM	86.81 $\pm$ 1.19	73.24 $\pm$ 1.55	93.98 $\pm$ 0.32	73.55 $\pm$ 0.22	79.81 $\pm$ 0.45	67.39 $\pm$ 1.73	53.11 $\pm$ 2.96	83.98 $\pm$ 0.97	78.16 $\pm$ 0.45	81.63 $\pm$ 0.46	77.17
Encoder	GCN <sub>LLMEmb</sub>	88.15 $\pm$ 1.79	76.45 $\pm$ 1.19	88.38 $\pm$ 0.68	74.39 $\pm$ 0.31	84.78 $\pm$ 0.86	68.27 $\pm$ 0.45	70.65 $\pm$ 0.75	84.23 $\pm$ 0.20	86.07 $\pm$ 0.20	89.52 $\pm$ 0.31	81.09
	ENGINE	87.00 $\pm$ 1.60	75.82 $\pm$ 1.52	90.08 $\pm$ 0.16	74.69 $\pm$ 0.36	85.44 $\pm$ 0.53	68.87 $\pm$ 0.25	71.21 $\pm$ 0.77	84.09 $\pm$ 0.09	86.98 $\pm$ 0.06	89.05 $\pm$ 0.13	81.32
Explainer	TAPE	88.05 $\pm$ 1.76	76.45 $\pm$ 1.60	93.00 $\pm$ 0.13	74.96 $\pm$ 0.14	87.11 $\pm$ 0.66	68.11 $\pm$ 0.54	66.22 $\pm$ 0.83	85.95 $\pm$ 0.59	87.72 $\pm$ 0.28	90.46 $\pm$ 0.18	81.80
Predictor	LLM <sub>IT</sub>	71.93 $\pm$ 1.47	60.97 $\pm$ 3.97	94.16 $\pm$ 0.19	76.08	80.61 $\pm$ 0.47	44.20 $\pm$ 3.06	58.30 $\pm$ 0.48	84.80 $\pm$ 0.13	78.27 $\pm$ 0.54	74.51 $\pm$ 0.53	72.38
	GraphGPT	82.29 $\pm$ 0.26	74.67 $\pm$ 1.15	93.54 $\pm$ 0.22	75.15 $\pm$ 0.14	82.54 $\pm$ 0.23	67.00 $\pm$ 1.22	60.72 $\pm$ 1.47	85.38 $\pm$ 0.72	84.46 $\pm$ 0.36	86.78 $\pm$ 1.14	79.25
	LLaGA	87.55 $\pm$ 1.15	76.73 $\pm$ 1.70	90.28 $\pm$ 0.91	74.49 $\pm$ 0.23	84.03 $\pm$ 1.10	69.16 $\pm$ 0.72	71.06 $\pm$ 0.38	85.56 $\pm$ 0.30	87.62 $\pm$ 0.30	90.41 $\pm$ 0.12	81.69

# Experiments & Findings

- Semi-supervised & Supervised

Semi-supervised		Cora	Citeseer	Pubmed	WikiCS	Instagram	Reddit	Books	Photo	Computer	Avg.
Classic	GCN <sub>ShallowEmb</sub>	82.30 $\pm$ 0.19	70.55 $\pm$ 0.32	78.94 $\pm$ 0.27	79.86 $\pm$ 0.19	63.50 $\pm$ 0.11	61.44 $\pm$ 0.38	68.79 $\pm$ 0.46	69.25 $\pm$ 0.81	71.44 $\pm$ 1.19	71.79
	SAGE <sub>ShallowEmb</sub>	82.27 $\pm$ 0.37	69.56 $\pm$ 0.43	77.88 $\pm$ 0.44	79.67 $\pm$ 0.25	63.57 $\pm$ 0.10	56.65 $\pm$ 0.33	72.01 $\pm$ 0.33	78.50 $\pm$ 0.15	81.43 $\pm$ 0.27	73.50
	GAT <sub>ShallowEmb</sub>	81.30 $\pm$ 0.67	69.94 $\pm$ 0.74	78.49 $\pm$ 0.70	79.99 $\pm$ 0.65	63.56 $\pm$ 0.04	60.60 $\pm$ 1.17	74.35 $\pm$ 0.35	80.40 $\pm$ 0.45	83.39 $\pm$ 0.22	74.67
	SenBERT-66M	66.66 $\pm$ 1.42	60.52 $\pm$ 1.62	36.04 $\pm$ 2.92	77.77 $\pm$ 0.75	59.00 $\pm$ 1.17	56.05 $\pm$ 0.41	83.68 $\pm$ 0.19	73.89 $\pm$ 0.31	70.76 $\pm$ 0.15	64.93
	RoBERTa-355M	72.24 $\pm$ 1.14	66.68 $\pm$ 2.03	42.32 $\pm$ 1.56	76.81 $\pm$ 1.04	63.52 $\pm$ 0.44	59.27 $\pm$ 0.34	84.62 $\pm$ 0.16	74.79 $\pm$ 1.13	72.31 $\pm$ 0.37	68.06
	GLEM	81.30 $\pm$ 0.88	68.80 $\pm$ 2.46	81.70 $\pm$ 1.07	76.43 $\pm$ 0.55	60.25 $\pm$ 3.66	55.13 $\pm$ 1.41	83.28 $\pm$ 0.39	76.93 $\pm$ 0.49	80.46 $\pm$ 1.45	73.81
Encoder	GCN <sub>LLMEmb</sub>	83.33 $\pm$ 0.75	71.39 $\pm$ 0.90	78.71 $\pm$ 0.45	80.94 $\pm$ 0.16	67.49 $\pm$ 0.43	68.65 $\pm$ 0.75	83.03 $\pm$ 0.34	84.84 $\pm$ 0.47	88.22 $\pm$ 0.16	78.51
	ENGINE	84.22 $\pm$ 0.46	72.14 $\pm$ 0.74	77.84 $\pm$ 0.27	80.94 $\pm$ 0.19	67.14 $\pm$ 0.46	69.67 $\pm$ 0.16	82.89 $\pm$ 0.14	84.33 $\pm$ 0.57	86.42 $\pm$ 0.23	78.40
Explainer	TAPE	84.04 $\pm$ 0.24	71.87 $\pm$ 0.35	78.61 $\pm$ 1.23	81.94 $\pm$ 0.16	66.07 $\pm$ 0.10	62.43 $\pm$ 0.47	84.92 $\pm$ 0.26	86.46 $\pm$ 0.12	89.52 $\pm$ 0.04	78.43
Predictor	LLM <sub>IT</sub>	67.00 $\pm$ 0.16	54.26 $\pm$ 0.22	80.99 $\pm$ 0.43	75.02 $\pm$ 0.16	41.83 $\pm$ 0.47	54.09 $\pm$ 1.02	80.92 $\pm$ 1.38	71.28 $\pm$ 1.81	66.99 $\pm$ 2.02	65.76
	GraphGPT	64.72 $\pm$ 1.50	64.58 $\pm$ 1.55	70.34 $\pm$ 2.27	70.71 $\pm$ 0.37	62.88 $\pm$ 2.14	58.25 $\pm$ 0.37	81.13 $\pm$ 1.52	77.48 $\pm$ 0.78	80.10 $\pm$ 0.76	70.02
	LLaGA	78.94 $\pm$ 1.14	62.61 $\pm$ 3.63	65.91 $\pm$ 2.09	76.47 $\pm$ 2.20	65.84 $\pm$ 0.72	70.10 $\pm$ 0.38	83.47 $\pm$ 0.45	84.44 $\pm$ 0.90	87.82 $\pm$ 0.53	75.07

Supervised		Cora	Citeseer	Pubmed	arXiv	WikiCS	Instagram	Reddit	Books	Photo	Computer	Avg.
Classic	GCN <sub>ShallowEmb</sub>	87.41 $\pm$ 2.08	75.74 $\pm$ 1.20	89.01 $\pm$ 0.59	71.39 $\pm$ 0.28	83.67 $\pm$ 0.63	63.94 $\pm$ 0.61	65.07 $\pm$ 0.38	76.94 $\pm$ 0.26	73.34 $\pm$ 1.34	77.16 $\pm$ 3.80	76.37
	SAGE <sub>ShallowEmb</sub>	87.44 $\pm$ 1.74	74.96 $\pm$ 1.20	90.47 $\pm$ 0.25	71.21 $\pm$ 0.18	84.86 $\pm$ 0.91	64.14 $\pm$ 0.47	61.52 $\pm$ 0.60	79.40 $\pm$ 0.45	84.59 $\pm$ 0.32	87.77 $\pm$ 0.34	78.64
	GAT <sub>ShallowEmb</sub>	86.68 $\pm$ 1.12	73.73 $\pm$ 0.94	88.25 $\pm$ 0.47	71.57 $\pm$ 0.25	83.94 $\pm$ 0.61	64.93 $\pm$ 0.75	64.16 $\pm$ 1.05	80.61 $\pm$ 0.49	84.84 $\pm$ 0.69	88.32 $\pm$ 0.24	78.70
	SenBERT-66M	79.61 $\pm$ 1.40	74.06 $\pm$ 1.26	94.47 $\pm$ 0.33	72.66 $\pm$ 0.24	86.51 $\pm$ 0.86	60.11 $\pm$ 0.93	58.70 $\pm$ 0.54	85.99 $\pm$ 0.58	77.72 $\pm$ 0.35	74.22 $\pm$ 0.21	76.40
	RoBERTa-355M	83.17 $\pm$ 0.84	75.90 $\pm$ 1.69	94.84 $\pm$ 0.06	74.12 $\pm$ 0.12	87.47 $\pm$ 0.83	63.75 $\pm$ 1.13	60.61 $\pm$ 0.24	86.65 $\pm$ 0.38	79.45 $\pm$ 0.37	75.76 $\pm$ 0.30	78.17
	GLEM	86.81 $\pm$ 1.19	73.24 $\pm$ 1.55	93.98 $\pm$ 0.32	73.55 $\pm$ 0.22	79.81 $\pm$ 0.45	67.39 $\pm$ 1.73	53.11 $\pm$ 2.96	83.98 $\pm$ 0.97	78.16 $\pm$ 0.45	81.63 $\pm$ 0.46	77.17
Encoder	GCN <sub>LLMEmb</sub>	88.15 $\pm$ 1.79	76.45 $\pm$ 1.19	88.38 $\pm$ 0.68	74.39 $\pm$ 0.31	84.78 $\pm$ 0.86	68.27 $\pm$ 0.45	70.65 $\pm$ 0.75	84.23 $\pm$ 0.20	86.07 $\pm$ 0.20	89.52 $\pm$ 0.31	81.09
	ENGINE	87.00 $\pm$ 1.60	75.82 $\pm$ 1.52	90.08 $\pm$ 0.16	74.69 $\pm$ 0.36	85.44 $\pm$ 0.53	68.87 $\pm$ 0.25	71.21 $\pm$ 0.77	84.09 $\pm$ 0.09	86.98 $\pm$ 0.06	89.05 $\pm$ 0.13	81.32
Explainer	TAPE	88.05 $\pm$ 1.76	76.45 $\pm$ 1.60	93.00 $\pm$ 0.13	74.96 $\pm$ 0.14	87.11 $\pm$ 0.66	68.11 $\pm$ 0.54	66.22 $\pm$ 0.83	85.95 $\pm$ 0.59	87.72 $\pm$ 0.28	90.46 $\pm$ 0.18	81.80
Predictor	LLM <sub>IT</sub>	71.93 $\pm$ 1.47	60.97 $\pm$ 3.97	94.16 $\pm$ 0.19	76.08	80.61 $\pm$ 0.47	44.20 $\pm$ 3.06	58.30 $\pm$ 0.48	84.80 $\pm$ 0.13	78.27 $\pm$ 0.54	74.51 $\pm$ 0.53	72.38
	GraphGPT	82.29 $\pm$ 0.26	74.67 $\pm$ 1.15	93.54 $\pm$ 0.22	75.15 $\pm$ 0.14	82.54 $\pm$ 0.23	67.00 $\pm$ 1.22	60.72 $\pm$ 1.47	85.38 $\pm$ 0.72	84.46 $\pm$ 0.36	86.78 $\pm$ 1.14	79.25
	LLaGA	87.55 $\pm$ 1.15	76.73 $\pm$ 1.70	90.28 $\pm$ 0.91	74.49 $\pm$ 0.23	84.03 $\pm$ 1.10	69.16 $\pm$ 0.72	71.06 $\pm$ 0.38	85.56 $\pm$ 0.30	87.62 $\pm$ 0.30	90.41 $\pm$ 0.12	81.69

*Encoder*  
*A robust choice*

*Explainer*  
*Suitable for graphs*  
*with labels heavily*  
*depend on text*

*Predictor*  
*Require rich*  
*supervision*

# Experiments & Findings

- Zero-shot

Type & LLM	Method	Cora (82.52)		WikiCS (68.67)		Instagram (63.35)		Photo (78.50)		Avg.	
		Acc	Macro-F1	Acc	Macro-F1	Acc	Macro-F1	Acc	Macro-F1	Acc	Macro-F1
LLM GPT-4o	Direct	68.08	69.25	68.59	63.21	44.53	42.77	63.99	61.09	61.30	59.08
	CoT	68.89	69.86	70.75	<b>66.23</b>	<b>47.87</b>	<b>47.57</b>	61.61	60.62	62.28	61.07
	ToT	68.29	69.13	70.78	65.69	44.16	42.68	60.84	59.16	61.02	59.16
	ReAct	68.21	69.28	69.45	66.03	44.49	43.16	63.63	60.82	61.44	59.82
	w. Neighbor	70.30	71.44	69.69	64.51	42.42	39.79	69.93	68.55	63.09	61.07
	w. Summary	<b>71.40</b>	<b>72.13</b>	<b>70.90</b>	65.42	45.02	44.62	<b>72.63</b>	<b>70.84</b>	<b>64.99</b>	<b>63.25</b>
LLM LLaMA-8B	Direct	62.64	63.02	56.77	53.04	37.58	29.70	41.23	44.26	49.56	47.50
	CoT	62.04	62.61	58.88	56.00	42.00	39.06	44.22	47.13	51.78	51.20
	ToT	34.06	33.30	40.35	41.15	<b>45.33</b>	<b>45.27</b>	31.31	34.00	37.76	38.43
	ReAct	36.55	38.04	22.40	25.76	44.67	44.42	27.03	28.96	32.66	34.30
	w. Neighbor	64.55	64.41	59.43	54.16	36.98	28.32	45.49	50.44	51.61	49.33
	w. Summary	<b>64.69</b>	<b>64.62</b>	<b>62.69</b>	<b>56.40</b>	37.59	30.91	<b>48.11</b>	<b>52.20</b>	<b>53.27</b>	<b>51.03</b>
GFM	ZeroG	<b>62.55</b>	<b>57.56</b>	<b>62.71</b>	<b>57.87</b>	<b>50.71</b>	<b>50.43</b>	46.27	<b>51.52</b>	<b>55.56</b>	<b>54.35</b>
	LLM <sub>IT</sub>	52.58	51.89	60.83	53.59	41.58	26.26	<b>49.23</b>	44.85	51.06	44.15
	LLaGA	18.82	8.49	8.20	8.29	47.93	47.70	39.18	4.71	28.53	17.30

*GFM*s can outperform open-source LLMs but still fall short of strong LLMs like GPT-4o

# Experiments & Findings

- LLM-as-Encoder vs. LM-as-Encoder

Method	Encoder	Semi-supervised				Supervised			
		Cornell	Texas	Wisconsin	Washington	Cornell	Texas	Wisconsin	Washington
Homophily Ratio (%)		11.55	6.69	16.27	17.07	11.55	6.69	16.27	17.07
MLP	SenBERT	50.59 $\pm$ 3.14	56.67 $\pm$ 2.15	71.98 $\pm$ 1.59	63.26 $\pm$ 2.89	66.15 $\pm$ 1.92	76.32 $\pm$ 3.72	81.51 $\pm$ 7.00	70.44 $\pm$ 8.65
	RoBERTa	59.08 $\pm$ 2.57	67.47 $\pm$ 1.29	73.87 $\pm$ 1.62	65.43 $\pm$ 3.44	66.67 $\pm$ 8.88	74.21 $\pm$ 6.09	80.00 $\pm$ 9.88	76.96 $\pm$ 7.48
	Qwen-3B	57.78 $\pm$ 3.24	76.27 $\pm$ 1.61	82.36 $\pm$ 1.62	75.11 $\pm$ 1.92	77.95 $\pm$ 4.76	88.95 $\pm$ 3.07	88.68 $\pm$ 6.64	83.48 $\pm$ 1.74
	Mistral-7B	59.87 $\pm$ 6.72	76.27 $\pm$ 1.08	83.30 $\pm$ 1.42	74.24 $\pm$ 0.88	78.46 $\pm$ 4.17	90.53 $\pm$ 3.16	89.43 $\pm$ 5.15	83.91 $\pm$ 5.60
GCN	SenBERT	46.80 $\pm$ 2.13	54.93 $\pm$ 0.68	58.30 $\pm$ 2.56	52.61 $\pm$ 1.35	50.77 $\pm$ 10.18	59.47 $\pm$ 5.16	61.13 $\pm$ 8.65	61.30 $\pm$ 1.62
	RoBERTa	47.06 $\pm$ 2.19	55.20 $\pm$ 2.78	54.91 $\pm$ 3.40	54.89 $\pm$ 1.50	51.79 $\pm$ 7.68	58.42 $\pm$ 7.33	59.24 $\pm$ 8.82	61.31 $\pm$ 5.39
	Qwen-3B	53.59 $\pm$ 2.07	56.80 $\pm$ 4.29	63.02 $\pm$ 2.16	64.56 $\pm$ 4.06	58.46 $\pm$ 10.56	64.74 $\pm$ 7.37	65.28 $\pm$ 6.82	67.83 $\pm$ 3.74
	Mistral-7B	54.64 $\pm$ 1.52	58.67 $\pm$ 3.60	62.08 $\pm$ 2.61	61.52 $\pm$ 3.61	59.49 $\pm$ 6.96	65.79 $\pm$ 6.66	64.90 $\pm$ 5.67	66.96 $\pm$ 4.84
SAGE	SenBERT	52.55 $\pm$ 1.58	61.73 $\pm$ 1.37	70.47 $\pm$ 1.75	65.54 $\pm$ 2.44	68.72 $\pm$ 4.97	80.00 $\pm$ 5.91	83.02 $\pm$ 6.31	76.96 $\pm$ 4.88
	RoBERTa	55.55 $\pm$ 3.44	64.26 $\pm$ 6.26	73.59 $\pm$ 2.72	66.08 $\pm$ 1.60	70.26 $\pm$ 8.37	80.53 $\pm$ 2.68	81.89 $\pm$ 7.42	74.35 $\pm$ 7.95
	Qwen-3B	57.13 $\pm$ 2.29	78.53 $\pm$ 1.76	83.21 $\pm$ 1.39	72.18 $\pm$ 3.66	74.87 $\pm$ 2.99	89.47 $\pm$ 1.67	91.32 $\pm$ 2.82	83.48 $\pm$ 3.25
	Mistral-7B	56.86 $\pm$ 1.37	76.53 $\pm$ 2.40	83.96 $\pm$ 1.55	73.91 $\pm$ 0.97	77.44 $\pm$ 2.99	91.05 $\pm$ 2.69	89.44 $\pm$ 4.24	81.74 $\pm$ 4.48
H <sub>2</sub> GCN	SenBERT	56.34 $\pm$ 1.67	66.67 $\pm$ 2.95	73.40 $\pm$ 1.68	70.55 $\pm$ 4.95	73.85 $\pm$ 7.14	84.21 $\pm$ 4.40	86.42 $\pm$ 6.01	77.83 $\pm$ 7.20
	RoBERTa	60.00 $\pm$ 3.54	68.13 $\pm$ 2.93	75.66 $\pm$ 2.12	71.52 $\pm$ 1.22	74.87 $\pm$ 7.68	83.16 $\pm$ 6.14	84.53 $\pm$ 9.04	79.13 $\pm$ 5.43
	Qwen-3B	61.57 $\pm$ 3.89	80.13 $\pm$ 6.45	84.53 $\pm$ 0.70	74.67 $\pm$ 1.77	76.41 $\pm$ 2.99	92.11 $\pm$ 2.88	89.81 $\pm$ 3.29	85.22 $\pm$ 3.99
	Mistral-7B	59.22 $\pm$ 4.54	72.93 $\pm$ 8.21	81.89 $\pm$ 1.51	68.59 $\pm$ 4.46	75.89 $\pm$ 3.84	89.47 $\pm$ 3.72	89.43 $\pm$ 5.42	86.09 $\pm$ 3.25

*LLM-as-Encoder significantly outperforms LMs in less informative graphs, e.g., **heterophilic ones***

# Contribution Summary

- A Testbed
  - LLMNodeBed, a PyG-based testbed including 14 datasets, 8 LLM-based algorithms, 8 classic algorithms, and 3 learning configurations
- Comprehensive Experiments
  - Training and evaluating over 2,700 models, we analyze how learning paradigm, homophily, language model type and size, and prompt design impact the performance
- Insights and Tips
  - Our work provides intuitive explanations, practical tips, and insights about the strengths and limitations of each algorithm category.

# Resources

- Paper: <https://arxiv.org/pdf/2502.00829>
- Code: <https://github.com/WxxShirley/LLMNodeBed>
- Dataset: <https://huggingface.co/datasets/xxwu/LLMNodeBed>
- Chinese Blog: <https://zhuanlan.zhihu.com/p/1913536056717967976>

Thank you for your attention!