Judge and Improve: Towards a Better Reasoning of Knowledge Graphs with Large Language Models

Anonymous ACL submission

Abstract

Graph Neural Networks (GNNs) have shown immense potential in improving the performance of large-scale models by effectively incorporating structured relational information. However, current approaches face two key challenges: (1) achieving robust semantic alignment between graph representations and large models, and (2) ensuring interpretability in the generated outputs. To address these challenges, we propose **ExGLM** (Explainable Graph Language Model), a novel training framework designed to seamlessly integrate graph and language modalities while enhancing transparency. Our framework introduces two core components: (1) a graph-language synergistic alignment module, which aligns graph structures with language model to ensure semantic consistency across modalities; and (2) a judge-andimprove paradigm, which allows the language model to iteratively evaluate, refine, and prioritize responses with higher interpretability, thereby improving both performance and transparency. Extensive experiments conducted on three benchmark datasets—ogbn-arxiv, Cora, and PubMed-demonstrate that ExGLM not only surpasses existing methods in efficiency but also generates outputs that are significantly more interpretable, effectively addressing the primary limitations of current approaches.

1 Introduction

002

007

017

021

042

Large Language Models (LLMs) have demonstrated remarkable success across various natural language processing tasks, including dialogue generation (Aboussalah and Ed-dib, 2025), machine translation (Zhu et al., 2024), question answering (Zhang et al., 2024), and text summarization (Zhang et al., 2025). These models exhibit an impressive capacity for understanding and generating human-like text. However, LLMs face inherent limitations in effectively modeling structured knowledge, such as graphs, which are essential



Figure 1: Overview of two mainstream methods (a) textualizing graph and inference via LLM and (b) aligning the semantic representation of LLMs and GNNs.

for capturing complex relationships and dependencies in diverse real-world domains like social networks, biological systems, and knowledge graphs. To address these limitations, recent research has explored the integration of GNNs (Kipf and Welling, 2017; Hamilton et al., 2018; Veličković et al., 2018) with LLMs (OpenAI et al., 2024; Yang et al., 2024; DeepSeek-AI et al., 2025), leveraging GNNs' strengths in modeling structured information alongside LLMs' powerful language capabilities, creating opportunities for enhanced performance in graph-related tasks.

Current approaches (Yang et al., 2021; Zhao et al., 2023; Xue et al., 2024) to combining GNNs and LLMs can be broadly classified into two categories. The first category involves **textualizing graph structures** and feeding them into LLMs (Figure 1(a)). For example, some methods (Zhao et al., 2023; Wang et al., 2024; Chen et al., 2024; Wu et al., 2025) describe nodes and their relationships using natural language templates to generate textual representations of subgraphs. Other approaches (Ye et al., 2024; Tang et al., 2024) employ special tokens to represent nodes and edges, effectively converting graph structures into sequences compatible with LLM processing. However, these methods have notable pitfalls: the textualization process can result in the loss of structural information, and the sequential representations may fail to fully capture the intricate relationships within the graph. Additionally, these methods face scalability challenges due to the token length constraints of LLMs, making them unsuitable for handling large graphs with extensive neighborhood information.

066

067

068

071

072

084

100

101

102

104

105

106

107

109

110

111

112 113

114

115

116

117

The second category of approaches (Chai et al., 2023; Tang et al., 2024; Liu et al., 2024) focuses on aligning the representation spaces of GNNs and LLMs in the semantic domain (Figure 1(b)). For instance, certain methods (Xia et al., 2024; Huang et al., 2023, 2024; Guo et al., 2025) project GNN-generated node embeddings into the embedding space of LLMs to achieve semantic consistency. Other techniques, such as those employing attention mechanisms (Ying et al., 2021; Kuang et al., 2022), integrate graph structure information directly into the language model's representations. While these approaches improve the integration of graph and language modalities, challenges remain. The alignment process may not be optimal, leading to performance bottlenecks in tasks requiring a precise understanding of graph structures and language semantics. Moreover, such methods often suffer from a lack of interpretability, making it difficult to elucidate how the model leverages graph information to make decisions or derive outputs.

To address the limitations of existing approaches, we propose ExGLM (Explainable Graph Language Model), a novel framework designed to effectively and interpretably integrate graph structures with LLMs. Our framework introduces a graph-language synergistic alignment module to achieve semantic consistency between graph structures and LLM outputs, while also maintaining interpretability. Specifically, we assign a textual attribute to each node in the graph, describing its adjacent relationships, with different nodes represented by special tokens. We then perform reasoning using the LLM and enhance its representation by incorporating the graph representation into the hidden state. To further improve interpretability, we propose a judge-and-improve paradigm where the LLM evaluates and selects responses with better interpretability. These optimized responses are subsequently used to refine the GNN-LLM model. Our main contribution can be summarized as follows:

• We propose a novel graph-language synergistic alignment module that effectively bridges the gap between graph-structured data and LLM outputs, ensuring robust semantic consistency across modalities. 118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

- We propose a judge-and-improve paradigm, enabling the model to iteratively evaluate and refine its responses for enhanced interpretability and generation quality, thereby improving both performance and transparency.
- We conduct comprehensive experiments on multiple datasets, demonstrating the superior performance and effectiveness of our approach compared to existing methods.

2 Related Work

2.1 Graph-Large Language Models

LLMs (OpenAI et al., 2024; Yang et al., 2024; DeepSeek-AI et al., 2025) achieve state-of-theart performance on various natural language tasks, however, it lacks explicit mechanisms to effectively model structured information, such as graphs. To address this limitation, recent studies (Shu et al., 2024; Tang et al., 2024) have explored ways to integrate the benefits of GNNs into LLM-based frameworks. For instance, Zhang et al. (2020) adapts the self-attention mechanism of BERT (Devlin et al., 2019) to capture the relational structure of nodes and edges within a graph. However, its performance is highly dependent on the presence and quality of node features, which may limit its applicability when such features are sparse or noisy. InstructGLM (Ye et al., 2024) leverages the natural language modeling capabilities of LLMs to describe multi-scale geometric structures within graphs, thereby improving representation and analysis of graph data. Nonetheless, it suffers from token-length limitations, making it challenging to process large graphs with extensive neighbor information. Jin et al. (2024) propose a framework named Graph-COT that enhances LLMs by encouraging them to perform iterative reasoning over graph structures. However, fine-tuning LLMs within this framework remains challenging, and potential misalignment between the graph structure and the text attribution can lead to inaccuracies. Another recent work, PromptGFM (Zhu et al., 2025),

explicitly prompts LLMs to mimic the workflow 166 of GNNs within the text space, achieving natu-167 rally alignment between graph representations and 168 textual modeling. While this approach improves 169 graph-text integration, it struggles to differentiate between graphs with similar semantic structures. 171 In this work, we propose a novel graph-language 172 synergistic alignment module that aligns GNNs 173 and LLMs at both the text attribution and semantic 174 representation levels. This alignment enables seam-175 less and effective incorporation of the strengths of 176 GNNs and LLMs. 177

2.2 Self-Judge-and-Improve Paradigm

178

The self-judge-and-improve paradigm highlights 179 the capacity of LLMs to autonomously evaluate 180 and enhance their own performance and capabili-181 ties, thereby reducing dependence on external su-182 pervision. This approach enables models to internally refine their understanding and outputs. For 185 instance, Self-Insturct (Wang et al., 2022) embodies this paradigm through a two-step process to 186 improve instruction-following abilities. First, the model generates sample outputs and evaluates them using its internal mechanisms, filtering out suboptimal results. These filtered samples are then 190 leveraged to fine-tune the model. Similarly, Self-191 Refine (Madaan et al., 2023) demonstrates how 192 LLMs can provide feedback on their own genera-193 tions and use this feedback to optimize their out-194 puts iteratively. Expanding on this concept, Yuan et al. (2025) introduced self-rewarding language 196 models, wherein LLMs assign self-generated rewards to their outputs. Preference pairs selected 198 based on these rewards are then utilized to opti-199 mize the models using DPO (Rafailov et al., 2023). While these approaches effectively minimize external intervention, the quality of self-judgment is inherently constrained by the performance of the LLM. To address this limitation, we propose the Judge-and-Improve paradigm, which incorporates a superior language model to evaluate the generated outputs. By introducing an external judgment 207 mechanism, our approach enhances the reliability and accuracy of evaluations, enabling more effec-210 tive refinement of the model's outputs.

3 Method

211

212The training framework of our method, illustrated213in Figure 2, is composed of two key modules:214(1) Graph-language synergistic alignment mod-

ule and (2) Judge-and-improve paradigm. The graph-language synergistic alignment module ensures effective integration between the GNNs and the LLMs by aligning textual attributes and semantic representations, thereby maintaining consistency across modalities. The judge-andimprove paradigm operates in two stages: first, it generates and selects accurate and explainable results through prompting, creating a supervised fine-tuning (SFT) and preference dataset; second, it uses these two datasets to optimize the model, progressively enhancing both performance and interpretability. 215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

3.1 Problem Setup

Graph structure. Generally, a graph can be formally defined as G = (V, E, X), where V = $\{v_1, v_2, \ldots, v_n\}$ represents the set of nodes, $E \subseteq$ $V \times V$ represents the set of edges, encoding pairwise relationships between nodes, and $X \in \mathbb{R}^{n \times d}$ is the node feature matrix. Each $x_i \in \mathbb{R}^d$ corresponds to the feature vector of node v_i , where d represents the dimensionality of the node features. Node classification with LLM. Consider a node classification problem over a graph G = (V, E, X), where the goal is to assign one of k discrete class labels to each node. Let $Y = \{1, 2, \dots, k\}$ denote the set of class labels. The training data consists of labeled examples (x_i, y_i) , where $x_i \in \mathbb{R}^d$ represents the graph feature vector of node $v_i \in V$, and $y_i \in Y$ is the corresponding class label. The objective is to learn a classifier $f: X \xrightarrow{G} Y$, such that $f(x_i)$ accurately predicts the class label y_i for each node. In this work, we first derive textual attributions T_v of each node v, capturing its structural and feature information in a textual format. We then leverage both the graph structure and an LLM to perform reasoning. Consequently, the classification objective is refined to learning a classifier $f: T \xrightarrow{G-LLM} Y$ where T represents the textual descriptions derived from the graph's structural and feature information. This approach integrates the representational strengths of both GNNs and LLMs, enabling a more interpretable and semantically rich node classification framework.

Classification with interpretability. In real-world applications where interpretability is paramount, it is essential for classification models to not only make accurate decisions but also provide clear explanations for those decisions. Therefore, our ultimate goal is to train a classifier that not only performs classification tasks but also generates an-



Graph-Language Synergistic Alignment Module

Figure 2: The training framework of ExGLM.

[Task introduction]

Classify the article according to its topic into one of the following categories: [Label Set]. Node represents academic paper with a specific topic, link represents a citation between the two papers. Pay attention to the multi-hop link relationship between the nodes.

[Node Info]

265

266

269

Example1: (<*node*_{*v*}>, *description*_{*v*}) is connected with [k-hop $\textbf{neighbor nodes} \verb|(<\!node_{n_1}\!\!>\!\!,\!description_{n_1}), ..., (<\!node_{n_l}\!\!>\!\!,\!description_{n_l})$ within k hops. Example2: (<nodev>, descriptionv) is connected with [k-hop neighbor nodes](< $node_{n_1}$ >, $description_{n_1}$), ..., (< $node_{n_i}$ >, $description_{n_i}$) within k hops through [nodes under the path, if k>1], respectively. [Instruct] **Example1**: Which category should ($< node_v >, description_v$) be

classified ast Example2: Should (<node_n>, title_n) be classified as [random label]?



alytical content to elucidate its decision-making process. This can be represented as a function $f: X \xrightarrow{G} Y$, Analysis. where Analysis provides the explanatory content.

3.2 **Graph-Language Synergistic Alignment** Module

To effectively leverage both structural information 271 from graphs and textual attributes from LLMs, we integrate GNNs and LLMs to obtain node repre-273 sentations. To bridge the gap between these two modalities, we propose a graph-language syner-275 gistic alignment module. This module consists of 276 two core components:(1) Textual attribution of adjacency relationships, which captures the textual representation of graph structures. (2) Incorporating graph semantic information into textual attribution, which enriches textual descriptions with 281 graph-based semantics. We detail these components below.

Textual attribution of adjacency relationships.

We derive the textual attribution of each node through a two-step process: (1) Subgraph sampling for node information. In the context of large graphs, subgraph sampling is crucial to mitigate computational complexity and enable scalable processing. In this work, we adopt a k-hop sampling strategy extract localized subgraphs centered around each node. Specifically, for a central node v, we sample its neighbors within k hops, and represent it as $\mathcal{N}_{v}^{(n)}$. v and $\mathcal{N}_{v}^{(n)}$ are then further utilized to derive the textual attribution of adjacency relationships.

289

290

291

292

294

295

296

297

300

301

302

303

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

(2) Node description via text. For each central node v, we construct multiple text descriptions. Each description is represented as a tuple T_v : (task introduction, <node info>, instruction). Here, task introduction provides a brief overview of the task, <node info> contains textual descriptions of the central node's neighbors sampled from *i*-hop neighborhoods (where $0 < i \leq k$, selected randomly), and instruction specifies a task-related question tailored to the node and its neighborhood information. Specifically, in the <node info> part, each individual graph node is represented by a special token $node_i$ with its brief text descriptions $description_i$. Detailed examples are provided in Figure 3. This approach enables the attribution of each graph node to be naturally expressed in textual form, bridging the structural information of graphs with the representational capabilities of LLMs.

Incorporating graph semantic information into textual attribution. Since the aforementioned special tokens for each node cannot be effectively represented by LLMs alone, we integrate them with representations derived from GNNs. To obtain

403

404

405

406

407

408

409

410

411

the representations of the GNNs, we adopt Graph-SAGE, which primarily consists of three steps: neighbor sampling, message aggregation, and node updating. For each node v, we sample its *m*-hop neighbors and denote them as $\mathcal{N}^{(m)}(v)$. In this paper, we set m = 2 in all scenarios. For the aggregation representation is calculated via mean-pooling of neighborhood features:

321

325

326

328

330

332

336

338

339

341

343

344

351

$$h_{\text{agg}}^{(l)} = \frac{1}{|\mathcal{N}^{(m)}(v)| + 1} \left(h_v^{(l-1)} + \sum_{u \in \mathcal{N}^{(m)}(v)} h_u^{(l-1)} \right)$$
(1)

in which $h_v^{(l-1)}$ denotes the representation of node v at layer (l-1) Then each node is updated via Nonlinear projection with learnable parameters which is denoted as:

$$h_v^{(l)} = \sigma \left(W^{(l)} \cdot h_{\text{agg}}^{(l)} \right), \qquad (2)$$

 $W^{(l)}$ is the layer-specific weight matrix, $sigma(\cdot)$ denotes the ReLU activation function, $d^{(l)}$ is the dimensionality at layer l.

After obtaining $h_v^{(l)}$, we directly add it to the LLM's hidden states corresponding to the special token v which is shown in Figure 2 left.

To achieve better alignment between the LLM and GNN in the semantic space, we perform joint training. First, we construct a dataset First $\mathcal{D}_{align} = \{(T_v, Y_v)\}, v \in V$, where $Y_v = [c_1, c_2, \ldots, c_n]$ denotes the label sequence associated with node v. The alignment is achieved using the Negative Log-Likelihood (NLL) loss function:

$$\mathcal{L}_{\text{align}} = -\sum_{t=1}^{n} \log P\left(c_t \mid c_{< t}, T_v; \theta_{\text{LLM}}, \theta_{\text{GNN}}\right),$$
(3)

where θ_{LLM} denotes the parameters of the LLM, and θ_{GNN} denotes the parameters of the GNN encoder.

3.3 Judge and Improve

Building upon the dual-projection constrained mechanism, we achieve a deep collaboration between Graph Neural Networks (GNNs) and Large Language Models (LLMs). Beyond mere decisionmaking, providing reasonable and trustworthy analyses significantly enhances the interpretability of these decisions, which is crucial for various realworld applications. To ensure the interpretability of model decisions, we require the LLMs to not only generate accurate answers but also provide comprehensive explanations for their decisions. However, we have observed that the explanations generated by the LLMs are often suboptimal, indicating a need for further training. Considering the challenge of obtaining training data with annotated explanations, we adopt a judge-and-improve paradigm (Yuan et al., 2025) to enhance the interpretability of the LLMs. Specifically, our approach involves the following steps:

(1) Generating multiple responses: For a given input, the LLM generates multiple responses, each 'accompanied by an explanation.

(2) Judging quality: Superior LLM acts as a judge to evaluate these responses, selecting the one that is not only accurate but also provides a reasonable explanation.

(3) Optimizing through annotated data: The generated responses and explanations are then used to optimized the LLM, thereby improving the quality of its explanations.

Responses generation. As illustrated in Figure 3, we construct the text attributes of node v using a tuple T_v (task introduction, <node info>, instruction). To assemble a high-quality and diverse dataset, we first replace the instruction with several predefined instruction templates that convey the same intent, denoted as T'_v . Subsequently, we generate a response Y'_v or each T'_v a set $\{(T'_v, Y'_v)\}$, for all $v \in V$.

Judging quality. We require the superior LLM such as GPT-4 to evaluate the generated responses based on three criteria: correctness of the response, adherence to instructions, and reasonableness of the explanation. If none of the samples meet all three criteria, we repeat the response generation procedure. Ultimately, for each T'_v , we obtain a set of candidates $(Y'_{v_{\text{best}}}, Y'_{v_1}, \ldots, Y'_{v_{n-1}})$, where *n* denotes the number of generated responses.

Optimizing through annotated data. Building upon the generated responses, we construct a supervised fine-tuning (SFT) dataset: $\mathcal{D}_{sft} = \{(T'_v, Y'_{v_{best}}\}, v \in V, \text{ which aims to teach the LLM} to learn the pattern of the best response. The loss$ is computed as follows:

$$\mathcal{L}_{\rm sft} = -\sum_{t=1}^{n} \log P\left(c'_t \mid c'_{< t}, T'_v; \theta_{\rm LLM}, \theta_{\rm GNN}\right),\tag{4}$$

where $Y'_{v_{\text{best}}} = [c'_{\text{best}_1}, \dots, c'_{\text{best}_n}]$ and best_i denotes the *i*-th token of the best response.

This procedure is trained alongside the alignment process, and the overall loss becomes:

$$\mathcal{L}_{\text{total}} = \lambda_1 \cdot \mathcal{L}_{\text{align}} + \lambda_2 \cdot \mathcal{L}_{\text{sft}}$$
(5)

where λ_1 and λ_2 denotes the hyperparameters.

412

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437 438

Moreover, to enable the LLM to distinguish be-413 tween good and bad responses, we construct a pref-414 erence dataset: $\mathcal{D}_{\text{pre}} = \{(T'_v, Y'_{v_{\text{pos}}}, Y'_{v_{\text{neg}}})\}, v \in V$ 415 where $Y'_{v_{\text{nos}}}$ denotes the best response (positive ex-416 ample), and $Y'_{v_{\text{new}}}$ denotes any other response (neg-417 ative example) corresponding to the same input T'_v . 418 This dataset pairs each best response with its cor-419 responding non-optimal responses for every node 420 $v \in V$. 421

We then utilize the preference dataset to optimize the LLM using DPO loss:

$$\mathcal{L}_{\text{DPO}} = -\mathbb{E}_{(T'_v, Y'_{v_{\text{pos}}}, Y'_{v_{\text{neg}}}) \sim \mathcal{D}_{\text{pre}}} \left[\log \sigma \left(\beta \log \frac{\pi_{\theta}(Y'_{v_{\text{pos}}} | T'_v)}{\pi_{\text{ref}}(Y'_{v_{\text{pos}}} | T'_v)} - \beta \log \frac{\pi_{\theta}((Y'_{v_{\text{neg}}} | T'_v)}{\pi_{\text{ref}}((Y'_{v_{\text{neg}}} | T'_v)} \right) \right],$$
(6)

where π_{ref} denotes the reference model, which we adopt as the model before DPO training, and β is a hyperparameter.

4 Experiments

4.1 Experimental Setup

Datasets. We utilize three graph datasets of varying scales: Cora, comprising 2,708 nodes and 5,429 edges (Yang et al., 2016); Pubmed, containing 19,717 nodes and 44,338 edges (Namata et al., 2012); and ogbn-arxiv, consisting of 169,343 nodes and 1,166,243 edges (Hu et al., 2020). For our experiments, we adopt the same dataset partitioning strategy as proposed in (Ye et al., 2024).

Metrics. Following (Namata et al., 2012), we use 439 accuracy as the primary metric to evaluate node 440 classification performance. To assess the inter-441 pretability of the generated outputs, we utilize GPT-442 4 (Brown et al., 2020) as an automated evaluator. 443 Additionally, to ensure a more robust and reliable 444 assessment of interpretability, we complement this 445 with a questionnaire-based survey, which provides 446 valuable human-centered insights (Sperrle et al., 447 2021). 448

Baselines. We compare the proposed method against three categories of existing approaches:
(1) GNN-based models, including GCN (Kipf and Welling, 2017),GraphSAGE(Hamilton et al., 2018),GAT(Veličković et al., 2018),TransGAT,
(Louis et al., 2020) etc.; (2) Transformer-based models, such as Graphormer(Ying et al., 2018)

2021),GT(Dwivedi and Bresson, 2021) and Coar-Former (Kuang et al., 2022); and (3) LLM-based models, such as InstructGLM (Ye et al., 2024). **Implementations**. In our implementation, we uti-

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

lize Llama-7B (Touvron et al., 2023) and Llama3.1-8B-Instruct (Touvron et al., 2023) as the LLM backbones. For Llama3.1-8B-Instruct, when it serves as the backbone for InstructGLM, we make minimal adjustments to the inputs to align with its requirements, such as embedding dialogue templates and mapping node IDs to token IDs. During the data generation and annotation phase, we use the Qwen2.5-7B-Instruct model (Bai et al., 2023) to generate decision-analysis content for node classification tasks. To ensure high-quality annotations, we further employ the more powerful Qwen2.5-72B-Instruct model (Bai et al., 2023) as a "super annotator," automatically refining and validating the generated analyses. The model training is performed on 8 A100 GPUs, and all experiments are conducted over 1-3 epochs.

4.2 Performance Comparison

Tables 1 compares the performance of various models on the Cora and PubMed datasets, showcasing the effectiveness of different approaches.

Accuracy on Cora dataset. Among the GNNbased methods, ACM-GCN+ achieves the best accuracy on the Cora dataset (89.75%). Transformersbased methods, on the other hand, generally exhibit relatively lower performance. Notably, the hybrid InstructGLM approach, which combines GNN and LLM techniques, is the most comparable to our method, achieving competitive performance with an accuracy of 87.08% on Cora. In contrast, our proposed method achieves 88.8% accuracy, surpassing all existing Transformers-based and GNN-LLM-based approaches.

Accuracy on PubMed dataset. On the PubMed dataset, InstructGLM sets a strong baseline with the best performance among prior methods. Our method outperforms all baselines, achieving a new state-of-the-art accuracy of 94.6%. These results highlight the superiority of our training framework. Accuracy on Ogbn-Arxiv dataset. Table 2 summarizes the performance of various models on the Ogbn-Arxiv dataset. Among traditional GNN-based approaches, DRGAT achieves the highest accuracy at 76.11%, outperforming simpler architectures such as GraphSAGE (74.35%) and GAT (74.15%), which exhibit moderate performance. Notably, methods that integrate large language

Method Cora (%) PubMed (%) Туре MixHop **GNN** 75.65 90.04 **GNN** 76.70 GAT 83.28 Geom-GCN **GNN** 85.27 90.05 SGC-v2 **GNN** 85.48 85.36 86.58 GraphSAGE GNN 86.85 GNN GCN 87.78 88.90 BernNet **GNN** 88.52 88 48 FAGCN **GNN** 88.85 89.98 **GCNII GNN** 88.93 89.80 RevGAT **GNN** 89.11 88.50 Snowball-V3 **GNN** 89.59 91.44 ACM-GCN+ 89.75 90.96 GNN Transformers 80.41 88.24 Graphormer Transformers 86.42 88.75 GT CoarFormer Transformers 88.69 89.75 InstructGLM GNN-LLM 87.08 93.84 **ExGLM GNN-LLM** 88.8 94.6

Table 1: Accuracy on Cora and PubMed datasets.

models (LLMs) with GNN frameworks surpass all conventional GNN models, demonstrating the po-508 tential of combining structured graph data with the 509 510 rich semantic understanding of LLMs. For instance, InstructGLM achieves an accuracy of 76.42%, further highlighting the effectiveness of this hybrid 512 approach. Our proposed method achieves the high-513 est overall accuracy at 77.4%, setting a new state-514 515 of-the-art performance on this task. This result underscores the advantages of our framework in 516 effectively leveraging both graph structures and 517 textual information to improve predictive perfor-518 mance. 519

Table 2: Accuracy on Ogbn-Arxiv dataset.

Method	Туре	Accuracy (%)
GAT	GNN	74.15
GraphSAGE	GNN	74.35
GCN	GNN	73.29
AGDN	GNN	76.02
RvGAT	GNN	75.90
DRGAT	GNN	76.11
InstructGLM	GNN-LLM	76.42
ExGLM	GNN-LLM	77.4

Accuracy with Other LLMs. Table 3 compares the performance of our proposed method against InstructGLM on the Cora and PubMed datasets, utilizing two different LLM backbones: LLaMA and LLaMA3. Two key observations can be drawn from the results: (1) Our method consistently outperforms the baseline InstructGLM across both datasets, regardless of the underlying LLM backbone. This demonstrates the robustness and ef-

Table 3: Performance comparison with different LLMs.

Method	Cora (%) PubMed (%)		
InstructGLM (LLaMA)	87.08	93.84	
Ours (LLaMA)	88.8	94.6	
InstructGLM (LLaMA3)	88.01	94.17	
ExGLM (LLaMA3)	89.30	94.42	

fectiveness of our approach. (2) The use of a more advanced backbone does not always guarantee a significant performance improvement. While both methods perform slightly better with LLaMA3 compared to LLaMA, the relative gain is marginal. Notably, when applying LLaMA3, the performance on the PubMed dataset drops slightly from 94.6% to 94.42%. This indicates that the integration mechanism and model design play a more critical role than simply using a stronger LLM.

GPT-4 evaluation. To evaluate the different methods more comprehensively, we use GPT-4 as a proxy for human judgment. Specifically, we task GPT-4 with performing pairwise evaluations to select the better response based on three key criteria: correctness of the response, adherence to instructions, and reasonableness of the explanation. The evaluation results, presented in Table 4, demonstrate that our method outperforms InstructGLM on both datasets. The low performance of Instruct-GLM may be attributed to its overfitting on the dataset, which can lead to less fluent or less adaptable language generation. Additionally, the integration of DPO enhances overall performance on both datasets.

Table 4: GPT-4 evaluation results with LLaMA3 as base model.

ExGLM vs.	Dataset Win (%) Lose (%)		
InstructGLM	Cora	81.61	0
	PubMed	92.45	0
ExGLM (w/o DPO)	Cora	8.46	6.80
	PubMed	5.81	4.42

4.3 Interpretability

7

DPO influence for accracy. In the judge-andimprove paradigm, DPO is utilized to prioritize generations that exhibit better interpretability. However, it remains essential to evaluate how this prioritization affects reasoning accuracy. The results presented in Table 5 demonstrate that enhancing interpretability does not compromise accuracy 549

550

551

552

553

554

555

556

557

558

559

560

561

529

530

520



Figure 4: A show case of explanation provided by (left) GNN-based method (2) ExGLM.

and may even lead to improvements in reasoning performance.

Table 5: Ablation study of DPO with LLaMA as base model.

Method	Cora (%)	PubMed (%)
ExGLM w/o dpo	88.92	94.85
ExGLM	90.03	94.75

A showcase. We present a showcase in Figure 4 to illustrate the interpretability of our method in comparison with GNN-based approaches. While GNN-based methods provide explanations for their reasoning through attention weights, these weights may not accurately capture the underlying inference process and can be challenging for humans to interpret. In contrast, our method generates natural language explanations directly, thereby enhancing comprehensibility and interpretability.

Human evaluation. We aim to evaluate whether the use of DPO in the judge-and-improve paradigm 575 enhances interpretability. However, assessing interpretability is challenging due to the lack of a 577 standardized metric. To address this, we conducted a human evaluation. Specifically, we designed a questionnaire involving 20 human participants, each answering 20 questions. Participants were asked to select the response they deemed more interpretable based on three key criteria: coherency, 584 logical consistency, and factuality. The results of this evaluation, presented in Table 6, demonstrate the effectiveness of our approach. The baseline InstructGLM suffers from overfitting on the training dataset, which harms its language generation capa-588

bilities and limits its ability to provide meaningful explanations.

589

591

592

593

594

595

596

597

598

600

601

602

603

604

605

606

607

608

Table 6: Human evaluation results with LLaMA3 as base model.

ExGLM vs.	Dataset	Win (%)	Lose (%)
InstructGLM	Cora PubMed	$100.00 \\ 100.00$	0 0
ExGLM (w/o DPO)	Cora PubMed	23.25 9.75	13.75 9.50

5 Conclusion

This work investigates how to better leverage Large Language Models (LLMs) for reasoning with structured data. Concretely, we aim to address two main limitations identified in recent studies: crossmodality alignment and interpretability. We propose a novel training framework named ExGLM, within which a graph-language synergistic alignment module is introduced to ensure semantic consistency across modalities. Additionally, we introduce a judge-and improve paradigm that adopts a superior language model to evaluate and select generated responses with better interpretability. The selected data is subsequently utilized to optimize the reasoning model. Experiments across various scenarios demonstrate the effectiveness of our approach, showcasing its potential to advance reasoning with structured data.

563

564

565

6 Limitations

609

627

631

632

633

634

635

636

637

638

641

643

644

647

649

654 655

659

While our work achieves promising results, there are several limitations that warrant attention. First, 611 the effectiveness of the judge-and-improve module depends heavily on the performance of the superior language model used for evaluation. If the 614 615 evaluating model introduces biases or provides inaccurate assessments, the refinement process may 616 be suboptimal, potentially constraining the overall 617 improvement of the target model's outputs. Second, the current framework does not implement the 619 620 judgment-and-improvement process iteratively. Iterative refinement, which involves multiple rounds of evaluation and optimization, could further enhance the quality and robustness of the model's 623 outputs. However, this remains an unexplored avenue and is left for future work. 625

References

- Amine Mohamed Aboussalah and Abdessalam Ed-dib. 2025. Are gnns doomed by the topology of their input graph? *Preprint*, arXiv:2502.17739.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, and 1 others. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Ziwei Chai, Tianjie Zhang, Liang Wu, Kaiqiao Han, Xiaohai Hu, Xuanwen Huang, and Yang Yang. 2023.
 Graphllm: Boosting graph reasoning ability of large language model. *Preprint*, arXiv:2310.05845.
- Runjin Chen, Tong Zhao, AJAY KUMAR JAISWAL, Neil Shah, and Zhangyang Wang. 2024. LLaGA: Large language and graph assistant. In *Forty-first International Conference on Machine Learning*.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, and 181 others. 2025. Deepseek-r1: Incentivizing reasoning capability in Ilms via reinforcement learning. *Preprint*, arXiv:2501.12948.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the*

North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers), pages 4171–4186. 660

661

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

- Vijay Prakash Dwivedi and Xavier Bresson. 2021. A generalization of transformer networks to graphs. *Preprint*, arXiv:2012.09699.
- Zirui Guo, Lianghao Xia, Yanhua Yu, Yuling Wang, Kangkang Lu, Zhiyong Huang, and Chao Huang. 2025. Graphedit: Large language models for graph structure learning. *Preprint*, arXiv:2402.15183.
- William L. Hamilton, Rex Ying, and Jure Leskovec. 2018. Inductive representation learning on large graphs. *Preprint*, arXiv:1706.02216.
- Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. 2020. Open graph benchmark: Datasets for machine learning on graphs. *Advances in neural information processing systems*, 33:22118–22133.
- Xuanwen Huang, Kaiqiao Han, Dezheng Bao, Quanjin Tao, Zhisheng Zhang, Yang Yang, and Qi Zhu. 2023. Prompt-based node feature extractor for fewshot learning on text-attributed graphs. *Preprint*, arXiv:2309.02848.
- Xuanwen Huang, Kaiqiao Han, Yang Yang, Dezheng Bao, Quanjin Tao, Ziwei Chai, and Qi Zhu. 2024. Can gnn be good adapter for llms? In *Proceedings* of the ACM Web Conference 2024, WWW '24, page 893–904, New York, NY, USA. Association for Computing Machinery.
- Bowen Jin, Chulin Xie, Jiawei Zhang, Kashob Kumar Roy, Yu Zhang, Zheng Li, Ruirui Li, Xianfeng Tang, Suhang Wang, Yu Meng, and Jiawei Han. 2024. Graph chain-of-thought: Augmenting large language models by reasoning on graphs. *Preprint*, arXiv:2404.07103.
- Thomas N. Kipf and Max Welling. 2017. Semisupervised classification with graph convolutional networks. *Preprint*, arXiv:1609.02907.
- Weirui Kuang, Zhen WANG, Yaliang Li, Zhewei Wei, and Bolin Ding. 2022. Coarformer: Transformer for large graph via graph coarsening.
- Zheyuan Liu, Xiaoxin He, Yijun Tian, and Nitesh V. Chawla. 2024. Can we soft prompt llms for graph learning tasks? In *Companion Proceedings of the ACM Web Conference 2024*, WWW '24, page 481–484, New York, NY, USA. Association for Computing Machinery.
- Steph-Yves Louis, Yong Zhao, Alireza Nasiri, Xiran Wang, Yuqi Song, Fei Liu, and Jianjun Hu. 2020. Graph convolutional neural networks with global attention for improved materials property prediction. *Physical Chemistry Chemical Physics*, 22(32):18141– 18148.

823

824

770

Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, and 1 others. 2023. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36:46534–46594.

713

714

715

719

725

726

727

733

734

735

736

737

738

739

740

741

742

743

745

746

747 748

749

750

751

757

758

759

760

761

764

- Gerald Namata, Ben London, Andrey Kolobov, German Mart'inez-Mu noz, and Kristian Kersting. 2012.
 Query-driven active surveying for collective classification. In *Proceedings of the 10th international conference on Knowledge discovery and data mining*, pages 446–461. Springer.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, and 262 others. 2024. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn.
 2023. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36:53728–53741.
- Dong Shu, Tianle Chen, Mingyu Jin, Chong Zhang, Mengnan Du, and Yongfeng Zhang. 2024. Knowledge graph large language model (kg-llm) for link prediction. *Preprint*, arXiv:2403.07311.
- Fabian Sperrle, Mennatallah El-Assady, Grace Guo, Rita Borgo, D Horng Chau, Alex Endert, and Daniel Keim. 2021. A survey of human-centered evaluations in human-centered machine learning. In *Computer Graphics Forum*, volume 40, pages 543–568. Wiley Online Library.
- Jiabin Tang, Yuhao Yang, Wei Wei, Lei Shi, Lixin Su, Suqi Cheng, Dawei Yin, and Chao Huang. 2024. Graphgpt: Graph instruction tuning for large language models. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '24, page 491–500, New York, NY, USA. Association for Computing Machinery.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, and 1 others. 2023. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph attention networks. *Preprint*, arXiv:1710.10903.
- Jianing Wang, Junda Wu, Yupeng Hou, Yao Liu, Ming Gao, and Julian McAuley. 2024. InstructGraph: Boosting large language models via graph-centric

instruction tuning and preference alignment. In *Find-ings of the Association for Computational Linguistics: ACL 2024*, pages 13492–13510, Bangkok, Thailand. Association for Computational Linguistics.

- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2022. Self-instruct: Aligning language models with self-generated instructions. *arXiv preprint arXiv:2212.10560*.
- Menghua Wu, Russell Littman, Jacob Levine, Lin Qiu, Tommaso Biancalani, David Richmond, and Jan-Christian Huetter. 2025. Contextualizing biological perturbation experiments through language. In *The Thirteenth International Conference on Learning Representations*.
- Lianghao Xia, Ben Kao, and Chao Huang. 2024. Open-Graph: Towards open graph foundation models. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 2365–2379, Miami, Florida, USA. Association for Computational Linguistics.
- Rui Xue, Xipeng Shen, Ruozhou Yu, and Xiaorui Liu. 2024. Efficient end-to-end language model fine-tuning on graphs. *Preprint*, arXiv:2312.04737.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, and 43 others. 2024. Qwen2 technical report. *Preprint*, arXiv:2407.10671.
- Junhan Yang, Zheng Liu, Shitao Xiao, Chaozhuo Li, Defu Lian, Sanjay Agrawal, Amit S, Guangzhong Sun, and Xing Xie. 2021. Graphformers: GNNnested transformers for representation learning on textual graph. In Advances in Neural Information Processing Systems.
- Zhilin Yang, William W. Cohen, and Ruslan Salakhutdinov. 2016. Revisiting semi-supervised learning with graph embeddings. *Preprint*, arXiv:1603.08861.
- Ruosong Ye, Caiqi Zhang, Runhui Wang, Shuyuan Xu, and Yongfeng Zhang. 2024. Language is all a graph needs. *EACL*.
- Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin Zheng, Guolin Ke, Di He, Yanming Shen, and Tie-Yan Liu. 2021. Do transformers really perform badly for graph representation? *Advances in neural information processing systems*, 34:28877–28888.
- Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Xian Li, Sainbayar Sukhbaatar, Jing Xu, and Jason Weston. 2025. Self-rewarding language models. *Preprint*, arXiv:2401.10020.
- Ce Zhang, Taixi Lu, Md Mohaiminul Islam, Ziyang Wang, Shoubin Yu, Mohit Bansal, and Gedas Bertasius. 2024. A simple llm framework for

825 long-range video question-answering. *Preprint*, arXiv:2312.17235.

827

828 829

830 831

832 833

834

835

836

837

838

839

840

841

842

843

844

845 846

847

848

849

- Jiawei Zhang, Haopeng Zhang, Congying Xia, and Li Sun. 2020. Graph-bert: Only attention is needed for learning graph representations. *Preprint*, arXiv:2001.05140.
- Yang Zhang, Hanlei Jin, Dan Meng, Jun Wang, and Jinghua Tan. 2025. A comprehensive survey on process-oriented automatic text summarization with exploration of llm-based methods. *Preprint*, arXiv:2403.02901.
- Jianan Zhao, Meng Qu, Chaozhuo Li, Hao Yan, Qian Liu, Rui Li, Xing Xie, and Jian Tang. 2023. Learning on large-scale text-attributed graphs via variational inference. In *The Eleventh International Conference* on Learning Representations.
- Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun Chen, and Lei Li. 2024. Multilingual machine translation with large language models: Empirical results and analysis. *Preprint*, arXiv:2304.04675.
- Xi Zhu, Haochen Xue, Ziwei Zhao, Wujiang Xu, Jingyuan Huang, Minghao Guo, Qifan Wang, Kaixiong Zhou, and Yongfeng Zhang. 2025. Llm as gnn: Graph vocabulary learning for textattributed graph foundation models. *arXiv preprint arXiv:2503.03313*.