LANGUAGE MODELS ARE GRAPH LEARNERS

Anonymous authors

Paper under double-blind review

Abstract

Language Models (LMs) are increasingly challenging the dominance of domainspecific models, including Graph Neural Networks (GNNs) and Graph Transformers (GTs), in graph learning tasks. Following this trend, we propose a novel approach that empowers off-the-shelf LMs to achieve performance comparable to state-of-the-art GNNs on node classification tasks, without requiring any architectural modification. By preserving the LM's original architecture, our approach retains a key benefit of LM instruction tuning: the ability to jointly train on diverse datasets, fostering greater flexibility and efficiency. To achieve this, we introduce two key augmentation strategies: (1) Enriching LMs' input using topological and semantic retrieval methods, which provide richer contextual information, and (2) guiding the LMs' classification process through a lightweight GNN classifier that effectively prunes class candidates. Our experiments on real-world datasets show that backbone Flan-T5 models equipped with these augmentation strategies outperform state-of-the-art text-output node classifiers and are comparable to topperforming vector-output node classifiers. By bridging the gap between specialized task-specific node classifiers and general LMs, this work paves the way for more versatile and widely applicable graph learning models. We will open-source the code upon publication.

1 INTRODUCTION

029 030 031

000

001 002 003

004

005 006 007

008 009

010

011

012

013

014

015

016

017

018

019

021

025

026 027 028

There is a growing trend of utilizing Language Models (LMs) for machine learning tasks across diverse domains. This approach has shown tremendous promise in areas such as computer vision (Desai & Johnson, 2021), audio processing (Mittal et al., 2021), and multimodal learning (Alayrac et al., 2022). In the domain of graph learning, recent efforts have begun to explore the capabilities of LMs in understanding and processing graph structures. Wang et al. (2023) showed that LMs can detect node connectivity and identify cycles within graphs, while Fatemi et al. (2024) explored LMs' ability to evaluate graph scale and identify connected components. Furthermore, InstructGLM (Ye et al., 2023) achieved state-of-the-art performance in text-output node classifiers for Text-Attributed Graphs (TAG) (Zhang et al., 2024a), where each node is associated with textual information.

040 However, InstructGLM suffers from a fundamental limitation that compromises the generality of the 041 backbone LM. Specifically, InstructGLM expands the LM's vocabulary by creating unique tokens 042 for each node, which incorporates topology-aware node embeddings as token embeddings. While 043 this approach is effective, it comes at the cost of reducing the LM's versatility, making it incompat-044 ible with two important use cases: (1) multi-task learning on diverse datasets, a common strategy for training Foundational Models (Wei et al., 2022; Chung et al., 2024), and (2) certain personalized LM fine-tuning services (Li et al., 2024b) that restrict modifications to the backbone model 046 architecture. This limitation raises a crucial question: How can off-the-shelf, text-to-text instruction-047 tuned LMs (Raffel et al., 2020) achieve competitive performance in node classification tasks without 048 requiring architectural modifications? 049

In stark contrast to (Huang et al., 2023), which suggests that LMs may only interpret graph structures in prompts as contextual paragraphs, our work presents a more optimistic outlook. We aim to overcome this inherent limitation by augmenting the input to LMs while preserving their original architecture. Our proposed model, named AUGGLM (Augmented Graph Language Model), leverages two key augmentation strategies to enhance the LM's ability to process graph data as follows:

- **Relevant Node Retrieval**: In contrast to InstructGLM, which relies on multi-hop ego networks akin to message-passing GNNs for structure-aware contextualization, AUG-GLM draws inspiration from Graph Transformers (GTs) (Min et al., 2022) and Retrieval-Augmented Generation (RAG) (Lewis et al., 2020; Guu et al., 2020). This enables the LM to access long-range structural and semantic information about the target node. We propose two complementary approaches to achieve this: (1) topological retrieval, and (2) prototypical semantic retrieval.
 - **Candidate Label Pruning**: To improve LMs' understanding of graph data while maintaining their text-to-text architecture, we convey the guidance from a specialist model, a pretrained lightweight GNN, to the input of LMs via narrowing down the potential labels. This allows LMs to focus on discerning between closely related candidates, ultimately enhancing the performance.

We conduct an extensive evaluation of our approach on four real-world TAGs, showing the effectiveness of our proposed instruction fine-tuning strategy. The results indicate that backbone LMs augmented with AUGGLM significantly outperform InstructGLM, while also matching or surpassing the performance of state-of-the-art vector-output classifiers. These findings represent a crucial step towards bridging the gap between specialized task-specific node classifiers and more general, fine-tuned LMs, highlighting the potential for unified models that can excel across multiple tasks.

073 074

054

056

058

060 061 062

063

064

065

066

067

069

071

072

075 076

2 RELATED WORK

077 078 079

LMs for Graphs. Recent studies have explored the ability of LMs to understand graph topology by investigating problems such as graph substructure recall (Wang et al., 2024), circle and connectivity 081 detection (Wang et al., 2023; Perozzi et al., 2024), node/edge counting (Perozzi et al., 2024), spatialtemporal problems on dynamic graphs (Zhang et al., 2024b). Notably, Fatemi et al. (2024) found 083 that the presentation format of graph data in text significantly impacts performance across various 084 tasks, highlighting the importance of effective graph-to-text encoding. Building on these findings, 085 several studies have explored classification tasks on TAGs, including node classification (Ye et al., 2023; Zhao et al., 2023b; Li et al., 2024a; Qin et al., 2023), link prediction (Brannon et al., 2023; 087 Tan et al., 2024), transfer learning (Tang et al., 2024), and graph reasoning (Jin et al., 2024); more-088 over, Chen et al. (2023) present a systematic summary on the performance of off-the-shelf solutions in two categories: LLMs-as-Enhancer and LLMs-as-Predictor. Furthermore, Zhang (2023) 089 proposed Graph-ToolFormer, a framework that enhances LMs with graph reasoning API tools, enabling more complex graph reasoning tasks. GIANT (Chien et al., 2022) and GLEM (Zhao et al., 091 2023a) utilize the interaction between graph data and LMs to generate better graph representations. 092 TAPE (He et al., 2024) leverages the explanations generated by an LLM to enrich the textual features, which are then used to fine-tune two LMs. The features from these LMs are subsequently 094 passed through an ensemble of GNNs for final prediction. 095

Retrieval-Augmented Generation (RAG). (Lewis et al., 2020; Karpukhin et al., 2020) enhance 096 LMs by granting them access to external knowledge (Hashimoto et al., 2018) during text generation. 097 This technique involves retrieving relevant documents from a large corpus and conditioning the LM 098 on both the input query and the retrieved information. Building on this, REALM Guu et al. (2020) pretrains the retriever and generator end-to-end. Subsequently, RETRO Borgeaud et al. (2022) ef-100 ficiently scales retrieval-enhanced autoregressive to large datasets. A crucial component of RAG's 101 success is the loss function proposed by Shi et al. (2023) which enables the retriever to be trained 102 with the LM viewed as a black box. HyDE Yu et al. (2023) uses hypothetical document generation 103 to improve retrieval in RAG systems. Furthermore, RAG has extended to multimodal settings Ya-104 sunaga et al. (2023). Recently, GraphRAG (Edge et al., 2024) has garnered significant attention 105 which involves constructing a Knowledge Graph (KG) and then generating responses based on the summaries of communities derived from KG. These advancements have significantly improved the 106 factual accuracy and contextual relevance of generated text, solidifying RAG as a promising tech-107 nique for various applications, including question answering and open-domain dialogue systems.

108 3 PRELIMINARIES

127 128

129

130

131

132 133

145

146

150 151

155

156 157

We use the following notation conventions: bold lower-case letters (e.g., x) denote column vectors, bold upper-case letters (e.g., X) denote matrices, and calligraphic upper-case letters (e.g., \mathcal{X}) denote sets. We use [·] and [·, ·] to index vectors and matrices, respectively.

113 We consider the node classification problem on TAGs where each node is associated with textual 114 attributes. A TAG with n nodes is formally represented as $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{T})$, where $\mathcal{V} = \{v_i\}_{i=1}^n$ 115 denotes a set of nodes, and $\mathcal{E} = \{e_{ij}\}_{i,j=1}^n$ is a set of edges where $e_{ij} = 1$ indicates that nodes v_i and 116 v_j are connected; otherwise, $e_{ij} = 0$. $\mathcal{T} = \{t_i\}_{i=1}^n$ indicates the set of node textual attributes. The 117 edges can also be represented by an adjacency matrix $\mathbf{A} \in \{0,1\}^{n \times n}$, where $\mathbf{A}[i,j] = 1$ if and only 118 if $e_{ij} = 1$. The training and test node labels are denoted by $\mathcal{Y} = \mathcal{Y}^{\text{train}} \cup \mathcal{Y}^{\text{test}} = \{y_i\}_{i=1}^n$, where 119 each label y_i belongs to one of the C classes, i.e., $y_i \in \{1, \ldots, C\}, \forall i$. In the semi-supervised setting 120 considered in this study, we have access to the graph structure and training labels $\mathcal{V}, \mathcal{E}, \mathcal{T}, \mathcal{Y}^{\text{train}}$ 121 during training. The task is to predict the labels of the remaining unlabeled nodes $\mathcal{Y}^{\text{test}}$.

Personalized PageRank (PPR) (Page, 1999; Jeh & Widom, 2003) assigns a relevance score to each node in the graph with respect to a given query node. Specifically, given the adjacency matrix A, the PPR scores $\mathbf{r}_i \in \mathbb{R}^n$ for all nodes in the graph relative to the query node v_i , are computed iteratively as follows:

$$\mathbf{r}_i \leftarrow (1 - \alpha) \tilde{\mathbf{A}} \mathbf{r}_i + \alpha \mathbf{q}_i \tag{1}$$

where $\alpha \in (0, 1)$ is the teleport probability, $\mathbf{q}_i \in (0, 1)^{n \times 1}$ is a one-hot vector with a single non-zero entry at index i, $\tilde{\mathbf{A}} = \mathbf{A}\mathbf{D}^{-1}$ is the normalized adjacency matrix, and \mathbf{D} is the degree matrix. Once the PPR scores are computed, we can identify the top-K most relevant nodes with respect to the query node v_i as follows:

$$PPR \operatorname{neigh}(v_i) = \{v_j : \mathbf{r}_i[j] \in \operatorname{top-K}(\mathbf{r}_i)\}$$
(2)

Language Models (LMs). We employ autoregressive LMs that predict the next token z_i based on the input sequence t and the context of previously generated tokens $z_{1:i-1}$. The probability of generating an N-token sequence z given the input sequence t is modeled as:

$$p_{\rm LM}(z|t) = \prod_{i=1}^{N} p_{\rm LM}(z_i|t, z_{1:i-1})$$
(3)

141 **Retrieval-Augmented Generation (RAG)** (Lewis et al., 2020; Guu et al., 2020) enhances the abil-142 ity of LMs to answer knowledge-intensive questions. Unlike traditional LMs that directly process 143 the input text t, RAG first retrieves a relevant document d^* from an external document corpus \mathcal{D} 144 using a similarity function s_{ϕ} :

$$d^* = \operatorname*{arg\,max}_{d \in \mathcal{D}} s_{\phi}(d, t) \tag{4}$$

147 The similarity function s_{ϕ} is typically implemented using a dual-encoder architecture (Bromley 148 et al., 1993) which computes the inner product $\langle \cdot, \cdot \rangle$ between the encoded representations of the 149 document and the input text:

$$s_{\phi}(d,t) = \langle \operatorname{Encoder}_{\phi}(d), \operatorname{Encoder}_{\phi}(t) \rangle$$
 (5)

Once the relevant document d^* is retrieved, it is fed into the LM along with the initial input tto estimate the output probability $p_{LM}(z|d^*, t)$. This approach enables LMs to leverage external knowledge and generate more accurate and informative responses.

4 Method

We explore the application of LMs to node classification tasks in graph learning, leveraging the instruction tuning paradigm to reformulate node classification as a text-to-text task Raffel et al. (2020).
Our method employs a carefully designed prompt template and a set of augmentation techniques to
transform graph data and ground truth labels into text pairs. This enables LMs to comprehend and
be fine-tuned for the task without requiring modifications to their underlying architecture.



Figure 1: Comparison of pipelines between the existing instruction-tuned LM InstructGLM and our approach, AUGGLM . Unlike InstructGLM, which explicitly encodes graph information into token embeddings as a form of soft prompting, AUGGLM maintains the original text-to-text framework of the off-the-shelf LM, offering greater generality and flexibility.



Figure 2: A detailed pipeline of AUGGLM (ours). In the semantic retrieval module, rectangle nodes denote the prototype of classes.

As shown in Figure 1, our approach differs from InstructGLM (Ye et al., 2023), the current stateof-the-art LM for TAGs, in its underlying design. While both methods utilize prompt templates to transform input graphs into text, InstructGLM relies on explicit encoding of node embeddings into the LM's token embeddings as a form of soft prompting (Lester et al., 2021). In contrast, our approach provides a more general framework, leveraging gradient descent through the LM without modifying its underlying text-to-text architecture. This design choice enables our model to retain the versatility of the original LM while adapting it to graph-based tasks. The following sections provide a detailed description of the specific techniques developed in the paper to achieve this goal.

199 200 201

202

189

190 191 192

193

194

195

196

197

169

170

171

172

4.1 RETRIEVAL-BASED AGGREGATION

General LMs are not designed to directly process graph-structured data. To overcome this limitation, a common approach is to employ prompt templates that transform graph data and associated tasks into a textual format that LMs can understand. For instance, consider the Cora (Sen et al., 2008) literature citation graph. A typical template (Huang et al., 2023; Ye et al., 2023) for node classification, as shown in Figure 3a consists of three main components: (1) a short description of the classification task, (2) the target node's textual features, such as its title and abstract, and (3) textual features from relevant neighboring nodes within the graph.

The success of the message-passing graph neural networks (GNNs) highlights the importance of the aggregation operation, whose typical examples are the sum and max pooling operations of intermediate node embeddings. A similar spirit is followed for the LM-based classifiers whose key design is at the selection of relevant nodes.

While existing works (Huang et al., 2023; Ye et al., 2023) select 1-hop or multi-hop neighbors as
 relevant nodes, we posit that this approach is suboptimal for two key reasons. Firstly, not all immediate or extended neighbors provide useful information for classifying the target node, which can

217

218

219

220 221

222

223

224

225 226

227

228 229

230 231

Model input Model input Please classify the following paper into Please classify the following paper based {pruned label candidates} based on the on the provided information. provided information. Title: {target node's title} Title: {target node's title} Content: {target node's abstract} Content: {target node's abstract} Related papers: {relevant nodes' titles} Related papers: {retrieved nodes' titles from PPR retrieval and/or semantic retrieval} Model output Model output {ground truth label} {ground truth label}

(a) A typical graph-to-text template

(b) Our template with augmented text features

Figure 3: Comparison of a typical graph-to-text template and our template with augmented text features.

236 introduce noise and degrade model performance. Secondly, incorporating multi-hop neighbors can 237 lead to "neighbor explosion" (Hamilton et al., 2017; Chen et al., 2018; Fey et al., 2021), i.e., an 238 exponentially-growing set of "relevant" nodes, resulting in increased computational costs, slower 239 model training/inference, and even out-of-memory issue. To address these limitations, we propose 240 two novel solutions for selecting relevant nodes: topological retrieval and prototypical semantic 241 retrieval. These methods are designed to efficiently identify the most informative nodes, thereby 242 enhancing the model's ability to capture meaningful graph structures and improve its overall perfor-243 mance.

Topological Retrieval. We leverage PPR (Page, 1999; Jeh & Widom, 2003) to perform topological retrieval, which has shown great effectiveness in conjunction with GNNs (Klicpera et al., 2019). The success of PPR suggests that the neighbors it identifies may provide more informative context than generic 1-hop or multi-hop neighbors. Specifically, for a target node v_i , we select its top-Kneighbors PPR neigh(v_i) based on their PPR scores, computed using Eqs. (1) and (2). The details of the PPR algorithm are introduced in Section 3. We then concatenate the text features from the PPR neighbors to form the PPR-retrieved text $t^{PPR retri} = \bigoplus_{j;v_j \in PPR neigh(v_i)} t_j$, where \oplus denotes text concatenation.

It is worth noting that the classic PPR algorithm is computationally expensive for large graphs due
to the matrix multiplication (Eq. (1)). However, efficient approximate solutions such as ApproximatePR (Andersen et al., 2006), can be applied to mitigate this issue. Nevertheless, PPR is a
topology-based heuristic that does not inherently leverage textual features nor adapt to feedback
from downstream node classification tasks. To address these limitations and enhance our framework's semantic awareness, we propose a complementary semantic retrieval strategy, which is discussed in the following section.

259 Prototypical Semantic Retrieval. Our semantic retrieval module draws inspiration from two pop-260 ular techniques: (1) RAG (Lewis et al., 2020; Guu et al., 2020), which leverage external corpora to 261 answer knowledge-intensive questions, and (2) Graph Transformers (Min et al., 2022) which aggregate messages from distant nodes via inner product-based attention weights. In the context of node 262 classification, we treat the textual features of all nodes except the target node as a surrogate "external 263 corpus." However, unlike typical question-answering tasks, retrieving textual features from a single 264 node is often insufficient for accurate node classification. To address this limitation, we enhance the 265 semantic retrieval process by focusing on prototypes, which capture the essence of each class (Snell 266 et al., 2017). 267

268 We employ prototypes (Biehl et al., 2016) as representative examples in the classification problem. 269 To obtain these prototypes, we first pretrain a lightweight GNN GNN_{ψ} , in a semi-supervised manner. This GNN generates a prediction vector for each node: $\tilde{\mathbf{y}}_i = \text{GNN}_{\psi}(v_i, \mathcal{G}) \in \mathbb{R}^c, \forall v_i \in \mathcal{V}$. We then compute the prediction confidence for each node v_i using the maximum logit:

$$\operatorname{Conf}(v_i) = \max_i \tilde{\mathbf{y}}_i[j] \in \mathbb{R}$$
(6)

The GNN's prediction for node v_i is denoted as $\tilde{y}_i = \arg \max_j \tilde{y}_i[j] \in \{1, \dots, C\}$, and the predicted class-*c* examples are $\tilde{\mathcal{Y}}_c = \{v_i : \tilde{y}_i = c\}$. We select the top-*N* confident examples as prototypes for each class *c*:

$$\mathcal{P}_{c} = \left\{ v_{i} : v_{i} \in \tilde{\mathcal{Y}}_{c} \land \operatorname{Conf}(v_{i}) \in \operatorname{top-N}\{\operatorname{Conf}(v_{j}) : v_{j} \in \tilde{\mathcal{Y}}_{c}\} \right\}$$
(7)

This process yields a total of $N \times C$ prototypes: $\mathcal{P} = \bigcup_{c \in \{1,...,C\}} \mathcal{P}_c$. To ensure that we retrieve text features from multiple nodes, we form our corpus \mathcal{D} by concatenating the text features of PPR neighbors for each prototype

$$\mathcal{D} = \left\{ \bigoplus_{j; v_j \in \mathsf{PPR neigh}(v_i)} t_j : v_i \in \mathcal{P} \right\}$$
(8)

Next, for each target node with its associated text features t^{target} , we compute the semantically retrieved text using Eq. (4): $t^{\text{semantic retri}} = \arg \max_{d \in \mathcal{D}} s_{\phi}(d, t^{\text{target}})$. In our experiments, we may use topological retrieval, prototypical semantic retrieval, or a hybrid approach that combines both by concatenating their retrieved texts. For simplicity, we denote the retrieved text as t^{retri} .

We defer the discussion of training ϕ and GNN_{ψ} to Section 4.4 and their specific architectures in Section 5.

4.2 CLASSIFIER GUIDANCE

Recent studies (Huang et al., 2023; Fatemi et al., 2024; Chen et al., 2024) highlighted the limited
understanding of graph topology in current mainstream LMs. While InstructGLM (Ye et al., 2023)
addresses this limitation by incorporating node embeddings from a pretrained GNN into the LM's
token embeddings, this approach necessitates modifications to the LM's architecture. We propose
an alternative method that conveys guidance from a pretrained GNN to inform LMs within the
input text space, thereby preserving the LM's original architecture while presenting a pruned set of
classification candidates.

We repurpose the pretrained GNN_{ψ} from the prototypical semantic retrieval module to inform AUG-GLM. For each node v_i , we identify and store the top-*I* predicted labels:

$$\mathcal{L}_{i} = \{j : \tilde{\mathbf{y}}_{i}[j] \in \text{top-I}(\tilde{\mathbf{y}}_{i})\} \in \{1, \dots, C\}^{I}$$
(9)

where I < C. Given that IndexToLabel maps are available for our target datasets, which associate numerical labels with their corresponding text representations, we can leverage this mapping to present the pruned label candidates for node v_i as a concatenated text: $t^{\text{candidates}} = \bigoplus_{i \in \mathcal{L}_i} \text{IndexToLabel}(i)$. The integration of this pruned candidate set into the input template is detailed in Section 4.3, where we elaborate on our overall template design.

By adopting this approach, we inject valuable structure-aware inductive bias from the GNN's prediction into the LM's input, thereby enhancing its ability to perform node classification tasks without
altering its fundamental architecture. By focusing the LM's attention on a smaller, more relevant set
of potential labels, we can improve its classification accuracy.

315 4.3 OVERALL TEMPLATE

Our augmented training samples are presented in Figure 3b, which includes three key elements: (1) the target node's text t^{target} , (2) the retrieved nodes' text t^{retri} , and (3) the pruned label candidates $t^{\text{candidates}}$. We collectively denote these elements as $t^{\text{input}} = (t^{\text{target}}, t^{\text{retri}}, t^{\text{candidates}})$. The backbone LM generates a prediction probability for the label sequence y^{target} according to the following equation:

314

303

304

272

277

278

$$p_{\rm LM}(y^{\rm target}|t^{\rm input}) = \prod_{i=1}^{|y|} p_{\rm LM}(y_i^{\rm target}|t^{\rm input}, y_{1:i-1}^{\rm target})$$
(10)

where $|\cdot|$ represents the sequence length; in this equation, *i* and 1: i - 1 are token-level indices. We will introduce the detailed selection of the backbone LM in Section 5. Figure 3b presents an exemplar template for the Cora dataset, showcasing the integration of t^{target} , t^{retri} , and $t^{\text{candidates}}$. A full list of templates used on all datasets in our experiments is detailed in Appendix C. Note that we exclude the abstracts of the retrieved nodes to prevent exceeding the maximum input length constraints of most LMs. During evaluation, we utilize only the "model input" portion of this template.

4.4 TRAINING

330 331

332

342 343

347 348 349

351 352 353

361 362

Our framework includes three parameterized modules that require training or fine-tuning: (1) GNNs for generating prototypes and candidate label pruning, as described in Sections 4.1 and 4.2, (2) the encoder ϕ from the semantic retriever, defined in Eq. 5, and (3) the backbone LM, utilized in Eq. 10. The GNNs from Sections 4.1 and 4.2 can be shared and their training is independent of the other modules which is supervised by ground truth labels. We provide more details on this process in Appendix A.

For the backbone LM, we adopt a standard training objective: minimizing the average token-wise negative log-likelihood (NLL) between the ground truth target sequence and the model's estimated output probability (Eq. 10). Specifically, for the target node, the NLL loss is computed as:

$$\mathcal{L}_{\rm NLL}(p_{\rm LM}(y^{\rm target}|t^{\rm input}), y^{\rm target})$$
(11)

To train the semantic retriever, we employ a distribution-matching loss. Specifically, for a given target node's text feature t^{target} , we first compute the retrieval probability distribution over all prototype text $t \in \mathcal{D}$:

$$p_{\phi}(t|t^{\text{target}}) = \frac{\exp(s_{\phi}(t, t^{\text{target}}))}{\sum_{t' \in \mathcal{D}} \exp(s_{\phi}(t', t^{\text{target}}))}$$
(12)

350 Next, we compute the empirical distribution supervised by the LM as:

$$\tilde{p}_{\rm LM}(t|t^{\rm target}, y^{\rm target}) = \frac{\exp(p_{\rm LM}(y^{\rm target}|t^{\rm target}, t))}{\sum_{t' \in \mathcal{D}} \exp(p_{\rm LM}(y^{\rm target}|t^{\rm target}, t'))}$$
(13)

This distribution represents the normalized importance of each prototype text $t \in D$ based on the LM's likelihood of generating the target ground truth label text y^{target} . We use \tilde{p} to distinguish this distribution from the generation probability defined in Eqs. (3) and (10). For simplicity, we omit the pruned candidate classes $t^{\text{candidates}}$ from the LM input in this equation, although they are indeed included in practice, as shown in Eq. (10).

The distribution matching loss is then computed as the Kullback-Leibler (KL) divergence between the retrieved distribution and the LM-supervised distribution:

KL (StopGradient
$$(\tilde{p}_{LM}(\cdot|t^{target}, y^{target})) || p_{\phi}(\cdot|t^{target}))$$
 (14)

This loss function aims to align the retrieved probability of each prototype text $t \in D$ with its importance in facilitating the LM's generation of the target label text y^{target} for the target node. The stop gradient operator ensures that the loss only updates the parameters of the semantic retriever ϕ , while keeping the LM's parameters θ frozen. This objective has been used by previous works (Shi et al., 2023; Izacard et al., 2023) without thorough analysis. We provide an in-depth examination of its properties and implications in Appendix B.

Notably, computing Eq. (13) requires $|\mathcal{D}|$ inferences of the backbone LM due to the denomi-369 nator. However, the LM is fine-tuned only on the NLL loss for the most relevant prototype, 370 $\arg \max_{d \in \mathcal{D}} s_{\phi}(d, t^{\text{target}})$ via Eq. (11). Consequently, each update step involves $|\mathcal{D}|$ forward passes 371 but only one backward pass. To further reduce the computational overhead associated with $|\mathcal{D}|$ in-372 ferences, we can employ a typical sampling strategy: selecting the top-M batch $\mathcal{D}_M = \{t : t \in$ 373 top- $M_{t' \in \mathcal{D}} s_{\phi}(t', t^{\text{target}})$ }. By replacing \mathcal{D} with \mathcal{D}_M in Eqs. (12) and (13), we can compute the re-374 trieval probability distribution and the LM-supervised distribution "in-batch", effectively reducing 375 the total number of inferences from $|\mathcal{D}|$ to M. 376

Algorithm 1 outlines a step-by-step process for fine-tuning our entire framework, processing one training node per step. This procedure can be readily extended to mini-batch settings.

1:	Given: (1) A graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{T})$ and training labels $\mathcal{Y}^{\text{train}}$, (2) initialized backbone LM θ , (3)
	initialized semantic encoder ϕ , and (4) initialized GNN ψ .
2:	Preprocessing : (1) train the GNN ψ based on $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{T})$ and $\mathcal{Y}^{\text{train}}$ till convergence; (2)
	generate prototypes and their text based on Eqs. (6), (7), and (8); (3) generate the pruned lab
	candidates for every node via Eq. (9).
3:	while θ and ϕ not converged do
4:	Sample $v_i \sim \mathcal{V}$.
5:	Retrieve the relevant nodes' text t_i^{retri} of node v_i via PPR retrieval (Eq. (2)) and/or semant
	retrieval (Eq. (4)).
6:	Plug t_i , t_i^{retri} , and $t_i^{\text{candidates}}$ (from preprocessing (3)) into the template (e.g., Figure 3b
	compute the NLL loss by Eq. (11), and update θ .
7:	Compute retrieval distribution $p_{\phi}(\cdot t_i)$ by Eq. (12).
8:	Call LM inference $ \mathcal{D} $ times to get $\{p_{\text{LM}}(y_i t_i,t)\}_{t\in\mathcal{D}}$ and $\tilde{p}_{\text{LM}}(\cdot t_i,y_i)$.
9:	Compute training loss for retriever by Eq. (14) and update ϕ .
10:	end while

397 Our model consists of three parameterized modules: (1) a GNN ψ for generating prototypes and 398 pruned label candidates, (2) the semantic retriever ϕ , and (3) the backbone LM θ . Notably, ψ 399 and ϕ are lightweight, with a number of parameters that is only 1/30 to 1/3 of the number of 400 parameters of LM θ . Compared to the state-of-the-art InstructGLM, our model has an additional 401 module ϕ , resulting in slightly more parameters which is relatively minor. During training, the GNN ψ can be trained independently, and the PPR scores can be precomputed. The training of θ 402 relies on the retrieved text from ϕ , while the training of ϕ requires $\tilde{p}_{LM}(\cdot|t^{\text{target}}, y^{\text{target}})$, which is 403 obtained through forward inference of θ . Importantly, their computational graphs (used for gradient 404 computation) are independent. This is because the LM θ concatenates the retrieved text from ϕ 405 into its input, which is not differentiable with respect to ϕ . Furthermore, when training ϕ , the loss 406 in Eq. (14) involves $\tilde{p}_{\rm LM}(\cdot|t^{\rm target}, y^{\rm target})$ wrapped with the StopGradient operator, ensuring that 407 the gradient computation does not update θ . As a result, the cost of back-propagation is similar to 408 updating the LM θ and the semantic encoder ϕ separately. 409

410

412

413

5 EXPERIMENTS

5.1 SETUP AND IMPLEMENTATION

Following (He et al., 2024; Ye et al., 2023), we evaluate our approach on four benchmark datasets: Cora (Sen et al., 2008), Pubmed (Sen et al., 2008), ogbn-arxiv (Hu et al., 2020), and a subset of ogbn-products (Hu et al., 2020; He et al., 2024). The statistics of the dataset are summarized in Table 5 (Appendix).

Our implementation employs two pretrained all-MiniLM-L6-v2 models¹ as the dual encoder for the 419 semantic retriever ϕ (Eq. (5)) and the text encoder for GNN ψ (Eq. (15)), respectively. We set the 420 teleport probability of the PPR to $\alpha = 0.1$. For the GNN, we employ a 3-layer GraphSAGE (Hamil-421 ton et al., 2017) architecture with a hidden dimension of 256 as ψ . Our hyperparameter settings 422 include K = 5 PPR neighbors, N = 10 prototypes, and M = 8 samples for LM inference. We 423 choose the number of label candidates I, from $\{2,3\}$. The backbone LM θ is implemented using 424 Flan-T5-small/base/large (Chung et al., 2022)², whose parameters are instruction-fine-tuned using 425 the templates shown in Figure 3b and Section C. 426

427 428

5.2 COMPARISON WITH STATE-OF-THE-ART

This section presents the comparison between AUGGLM and the state-of-the-art baselines.

¹https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2

²https://huggingface.co/docs/transformers/en/model_doc/flan-t5

	Method	Cora	Pubmed	ogbn-arxiv	ogbn-products
	GCN	87.78±0.96	88.90±0.32	73.60±0.18	75.64±0.21
	GraphSAGE	86.51±2.36	$89.08 {\pm} 0.28$	$73.88 {\pm} 0.33$	$76.04 {\pm} 0.25$
	BernNet	$88.52 {\pm} 0.95$	$88.48 {\pm} 0.41$	_	_
out	FAGCN	88.85±1.36	$89.98 {\pm} 0.52$	_	_
ut	GCNII	88.98±1.33	$89.80 {\pm} 0.52$	$72.74{\pm}0.16$	_
1-C	ACM-GCN	89.75±1.16	$91.44{\pm}0.59$	_	_
cto	GLEM + RevGAT	$88.56 {\pm} 0.60$	$94.71 {\pm} 0.20$	$76.97 {\pm} 0.19$	_
Ve	GIANT + RevGAT	83.53±0.38	$85.02 {\pm} 0.48$	$75.90{\pm}0.19$	$71.89 {\pm} 0.30$
	GIANT + GCN	84.23±0.53	$84.19 {\pm} 0.50$	$73.29 {\pm} 0.10$	$69.77 {\pm} 0.42$
	DeBERTa	76.06 ± 3.78	$94.94{\pm}0.46$	$73.61 {\pm} 0.04$	$72.97 {\pm} 0.23$
	TAPE + RevGAT	92.90±3.07	96.18±0.53	77.50±0.12	82.34±0.36
nt	ChatGPT-3.5	67.90	93.42	73.40	74.40
ıtpı	InstructGLM	90.77±0.52	$94.62 {\pm} 0.13$	$75.70{\pm}0.12$	_
-o	AUGGLM (T5-small)	91.14±0.55	$94.80{\pm}0.15$	$75.39{\pm}0.21$	$81.73 {\pm} 0.08$
ext	AUGGLM (T5-base)	91.24±0.46	$95.03 {\pm} 0.35$	76.80±0.14	81.91±0.11
Ľ	AUGGLM (T5-large)	91.51±0.26	95.16±0.18	$76.00{\pm}0.23$	82.90±0.10

Table 1: Performance comparison (accuracy) between AUGGLM and state-of-the-art models. The best-performing vector-output and text-output models are highlighted in **blue** and **red**, respectively.

We categorize models into two groups: (1) vector-output models which output a vector with di-454 mension equal to the number of classes, and (2) text-output models, which generate text as their 455 output. Specifically, we report results from GCN (Kipf & Welling, 2017), BernNet (He et al., 456 2021a), FAGCN (Bo et al., 2021), GCNII (Chen et al., 2020), ACM-GCN (Luan et al., 2022), 457 and GLEM (Zhao et al., 2023a)+RevGAT from the leaderboards³⁴⁵ and their published papers. The 458 results for TAPE+RevGAT, GIANT (Chien et al., 2022)+RevGAT (Li et al., 2021), GIANT+GCN, 459 DeBERTa (He et al., 2021b), and ChatGPT3.5 are reported from (He et al., 2024). The results of 460 InstructGLM are reported from (Ye et al., 2023). Note that all models, except ChatGPT-3.5, are 461 fine-tuned on the training set. We report mean and standard deviation over five runs. For text-output 462 models, we evaluate accuracy by checking whether the model's generated text matches the ground 463 truth label text exactly.

464 Table 1 presents a comprehensive comparison of our model's performance against existing ap-465 proaches. The results demonstrate that our proposed method consistently outperforms InstructGLM, 466 achieving new state-of-the-art performance among LMs with text-space outputs for node classifica-467 tion tasks on TAGs. Notably, this superior performance is achieved without modifying the under-468 lying architecture of the LMs, demonstrating the effectiveness of our approach. Furthermore, our 469 models exhibit competitive performance compared to the best vector-output models. Specifically, 470 on Cora, Pubmed, and ogbn-arxiv datasets, our models' performance closely approaches that of the state-of-the-art vector-output models. Furthermore, on the ogbn-products dataset, our approach 471 surpasses the performance of the best vector-output model, TAPE. 472

473 474

475

5.3 ABLATION STUDY

To evaluate the contribution of each key component in AUGGLM, we conducted an ablation study on three crucial modules: (1) topological retrieval, (2) semantic retrieval, and (3) candidate label pruning. We use the Flan-T5-small as the backbone LM for this analysis. The results, presented in Table 2, demonstrate that each module consistently improves performance across all datasets. Notably, our analysis reveals that the relative importance of each component varies across different datasets. For instance, candidate label pruning has a significant impact on performance for the Cora dataset, whereas its effect is less pronounced for the ogbn-products dataset. This variation in com-

483

485

9

432

433

434 435 436

³https://paperswithcode.com/sota/node-classification-on-cora-60-20-20-random

⁴https://paperswithcode.com/sota/node-classification-on-pubmed-60-20-20-random ⁵https://ogb.stanford.edu/docs/leader_nodeprop/

Table 2: Ablation study results. T, S, L, denotes the topological retrieval, semantic retrieval, and candidate label pruning. The \downarrow symbol denotes the decrease in accuracy of the ablated version compared to the full model.

	Т	S	L	Cora	Pubmed	ogbn-arxiv	ogbn-products
-	√ √	√	√	85.52 (↓5.62) 87.27 (↓3.87)	94.40 (\downarrow 0.40) 94.32 (\downarrow 0.48) 94.26 (\downarrow 0.54)	72.91 (↓2.48) 73.79 (↓1.60) 72.46 (↓1.03)	79.83 (\downarrow 1.90) 81.05 (\downarrow 0.68) 70.06 (\downarrow 2.67)
	\checkmark	\checkmark	\checkmark	90.23 (±0.89) 91.14	94.20 (↓0.34) 94.80	75.39	79.00 (<i>↓</i> 2.07) 81.73

Table 3: Performance of the jointly trained model (accuracy).

Training	Cora	Pubmed	ogbn-arxiv	ogbn-products
Multi-Task	91.52	94.52	74.87	82.29
Independent	91.14	94.80	75.39	81.73

ponent importance underscores adaptability of our approach, which can effectively accommodate diverse datasets with different characteristics.

5.4 MULTI-TASK TRAINING

One of the key advantage of pure text-to-text instruction tuning is that a single model can be trained on multiple tasks with the same input-output format. To verify this, we conducted an experiment using a Flan-T5-small model, applying our proposed strategies to jointly train it on four diverse datasets: Cora, Pubmed, ogbn-arxiv, and ogbn-products. The results, presented in Table 3 show that the jointly trained model achieve performance comparable to models trained separately on each individual dataset. We observe that on some datasets, such as Cora and ogbn-products, the jointly trained model even outperforms its dataset-specific counterparts.

These findings suggest that our approach can effectively handle multiple graph datasets using a single model, without incurring significant performance losses compared to models trained individually. This capability is crucial for efficient model deployment when dealing with diverse graph data. In contrast, other approaches, such as InstructGLM, require the addition of a large token dictionary to accommodate all nodes in the joint dataset, which hinders their ability to achieve similar generality. Moreover, most vector-output models, including TAPE, are limited by their predefined input-output dimensions, making them inflexible and unable to handle multiple datasets.

522 523

524

504

505 506

507

6 CONCLUSION

525 We introduce a novel framework called AUGGLM for instruct-tuning Language Models (LMs) to 526 perform node classification tasks on Text-Attributed Graphs (TAGs) using pure text-to-text instructions. Our approach is built upon two key innovations: (1) topological and semantic retrieval of 527 relevant nodes, (2) using a lightweight GNN to guide the classification process of LMs. Extensive 528 experimental results demonstrated the effectiveness of our framework, which consistently outper-529 formed the best existing text-output node classifiers, while achieving performance comparable to 530 state-of-the-art vector-output node classifiers. These findings suggest a promising direction for har-531 nessing the power of Large Language Models (LLMs) in graph learning tasks. 532

Limitation and Future Work. One limitation of this work is the need for manual definition of prompt templates in Table 4. A promising direction for future research is to develop methods for automatically searching for optimal templates in a data-driven manner. Another limitation is the requirement for pretraining a GNN ψ on each dataset, which stems from the inherent challenges of language models in understanding graph data. Addressing this limitation by developing more powerful language models capable of handling graph data is a challenging yet impactful area of future work.

540 REFERENCES

552

- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, 542 Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza 543 Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Mon-544 teiro, Jacob L. Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén 546 Simonyan. Flamingo: a visual language model for few-shot learning. In Sanmi Koyejo, 547 S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), Advances in Neu-548 ral Information Processing Systems 35: Annual Conference on Neural Information Process-549 ing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 550 2022, 2022. URL http://papers.nips.cc/paper files/paper/2022/hash/ 551 960a172bc7fbf0177ccccbb411a7d800-Abstract-Conference.html.
- Reid Andersen, Fan R. K. Chung, and Kevin J. Lang. Local graph partitioning using pagerank vectors. In 47th Annual IEEE Symposium on Foundations of Computer Science (FOCS 2006), 21-24 October 2006, Berkeley, California, USA, Proceedings, pp. 475–486. IEEE Computer Society, 2006. doi: 10.1109/FOCS.2006.44. URL https://doi.org/10.1109/FOCS.2006.44.
- Michael Biehl, Barbara Hammer, and Thomas Villmann. Prototype-based models in machine learn *Wiley Interdisciplinary Reviews: Cognitive Science*, 7(2):92–111, 2016.
- Deyu Bo, Xiao Wang, Chuan Shi, and Huawei Shen. Beyond low-frequency information in graph convolutional networks. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*, pp. 3950–3957. AAAI Press, 2021. doi: 10.1609/AAAI.V35I5.16514.
 URL https://doi.org/10.1609/aaai.v35i5.16514.
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Mil-566 lican, George van den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, Diego 567 de Las Casas, Aurelia Guy, Jacob Menick, Roman Ring, Tom Hennigan, Saffron Huang, Loren 568 Maggiore, Chris Jones, Albin Cassirer, Andy Brock, Michela Paganini, Geoffrey Irving, Oriol 569 Vinyals, Simon Osindero, Karen Simonyan, Jack W. Rae, Erich Elsen, and Laurent Sifre. Im-570 proving language models by retrieving from trillions of tokens. In Kamalika Chaudhuri, Ste-571 fanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (eds.), International 572 Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA, 573 volume 162 of Proceedings of Machine Learning Research, pp. 2206–2240. PMLR, 2022. URL 574 https://proceedings.mlr.press/v162/borgeaud22a.html. 575
- William Brannon, Suyash Fulay, Hang Jiang, Wonjune Kang, Brandon Roy, Jad Kabbara, and Deb
 Roy. Congrat: Self-supervised contrastive pretraining for joint graph and text embeddings. *CoRR*, abs/2305.14321, 2023. doi: 10.48550/ARXIV.2305.14321. URL https://doi.org/10.48550/arXiv.2305.14321.
- Jane Bromley, James W. Bentz, Léon Bottou, Isabelle Guyon, Yann LeCun, Cliff Moore, Eduard Säckinger, and Roopak Shah. Signature verification using A "siamese" time delay neural network. *Int. J. Pattern Recognit. Artif. Intell.*, 7(4):669–688, 1993. doi: 10.1142/S0218001493000339.
 URL https://doi.org/10.1142/S0218001493000339.
- Jianfei Chen, Jun Zhu, and Le Song. Stochastic training of graph convolutional networks with
 variance reduction. In Jennifer G. Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pp. 941–949. PMLR, 2018. URL http://proceedings.mlr.press/v80/chen18p.html.
- Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. Simple and deep graph convolutional networks. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pp. 1725–1735. PMLR, 2020. URL http://proceedings.mlr.press/v119/chen20v.html.

- Nuo Chen, Yuhan Li, Jianheng Tang, and Jia Li. Graphwiz: An instruction-following language model for graph computational problems. In Ricardo Baeza-Yates and Francesco Bonchi (eds.), *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, *KDD 2024, Barcelona, Spain, August 25-29, 2024*, pp. 353–364. ACM, 2024. doi: 10.1145/ 3637528.3672010. URL https://doi.org/10.1145/3637528.3672010.
- Zhikai Chen, Haitao Mao, Hang Li, Wei Jin, Hongzhi Wen, Xiaochi Wei, Shuaiqiang Wang, Dawei Yin, Wenqi Fan, Hui Liu, and Jiliang Tang. Exploring the potential of large language models (Ilms)in learning on graphs. *SIGKDD Explor.*, 25(2):42–61, 2023. doi: 10.1145/3655103.3655110.
 URL https://doi.org/10.1145/3655103.3655110.
- Eli Chien, Wei-Cheng Chang, Cho-Jui Hsieh, Hsiang-Fu Yu, Jiong Zhang, Olgica Milenkovic, and Inderjit S. Dhillon. Node feature extraction by self-supervised multi-scale neighborhood prediction. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. URL https://openreview.net/ forum?id=KJggliHbs8.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Y. Zhao, Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language models. *CoRR*, abs/2210.11416, 2022. doi: 10.48550/ARXIV.2210.11416.
- 616 Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan 617 Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, 618 Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pel-619 lat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Y. Zhao, 620 Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, 621 Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language 622 models. J. Mach. Learn. Res., 25:70:1-70:53, 2024. URL https://jmlr.org/papers/ v25/23-0870.html. 623
- 624 Karan Desai and Justin Johnson. Virtex: Learning visual representations from textual an-625 In IEEE Conference on Computer Vision and Pattern Recognition, CVPR notations. 626 2021, virtual, June 19-25, 2021, pp. 11162-11173. Computer Vision Foundation / 627 IEEE, 2021. doi: 10.1109/CVPR46437.2021.01101. URL https://openaccess. 628 thecvf.com/content/CVPR2021/html/Desai_VirTex_Learning_Visual_ 629 Representations_From_Textual_Annotations_CVPR_2021_paper.html.
- Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt,
 and Jonathan Larson. From local to global: A graph rag approach to query-focused summariza-*arXiv preprint arXiv:2404.16130*, 2024.

634

635

636

637

- Bahare Fatemi, Jonathan Halcrow, and Bryan Perozzi. Talk like a graph: Encoding graphs for large language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.* OpenReview.net, 2024. URL https://openreview.net/forum?id=IuXR1CCrSi.
- Matthias Fey, Jan Eric Lenssen, Frank Weichert, and Jure Leskovec. Gnnautoscale: Scalable and expressive graph neural networks via historical embeddings. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 3294–3304. PMLR, 2021. URL http://proceedings.mlr.press/v139/fey21a.html.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. Retrieval aug mented language model pre-training. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pp. 3929–3938. PMLR, 2020. URL http://proceedings.mlr.press/v119/guu20a.html.

669

- 648 William L. Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning 649 In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. on large graphs. 650 Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), Advances in Neural Information Processing Systems 30: Annual Conference on Neural Infor-mation Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pp. 651 652 1024-1034, 2017. URL https://proceedings.neurips.cc/paper/2017/hash/ 653 5dd9db5e033da9c6fb5ba83c7a7ebea9-Abstract.html. 654
- 655 Tatsunori B. Hashimoto, Kelvin Guu, Yonatan Oren, and Percy Liang. A retrieve-and-656 edit framework for predicting structured outputs. In Samy Bengio, Hanna M. Wallach, 657 Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett (eds.), Ad-658 vances in Neural Information Processing Systems 31: Annual Conference on Neural Infor-659 mation Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada, 660 pp. 10073-10083, 2018. URL https://proceedings.neurips.cc/paper/2018/ 661 hash/cd17d3ce3b64f227987cd92cd701cc58-Abstract.html.
- Mingguo He, Zhewei Wei, Zengfeng Huang, and Hongteng Xu. Bernnet: Learning arbi-663 trary graph spectral filters via bernstein approximation. In Marc'Aurelio Ranzato, Alina 664 Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), Ad-665 vances in Neural Information Processing Systems 34: Annual Conference on Neural Infor-666 mation Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual, pp. 14239-667 14251, 2021a. URL https://proceedings.neurips.cc/paper/2021/hash/ 668 76f1cfd7754a6e4fc3281bcccb3d0902-Abstract.html.
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. Deberta: decoding-enhanced 670 bert with disentangled attention. In 9th International Conference on Learning Representations, 671 ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net, 2021b. URL https: 672 //openreview.net/forum?id=XPZIaotutsD. 673
- 674 Xiaoxin He, Xavier Bresson, Thomas Laurent, Adam Perold, Yann LeCun, and Bryan Hooi. Har-675 nessing explanations: Llm-to-lm interpreter for enhanced text-attributed graph representation 676 learning. In The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024. OpenReview.net, 2024. URL https://openreview.net/ 677 forum?id=RXFVcynVe1. 678
- 679 Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, 680 and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs. CoRR, 681 abs/2005.00687, 2020. URL https://arxiv.org/abs/2005.00687. 682
- 683 Jin Huang, Xingjian Zhang, Qiaozhu Mei, and Jiaqi Ma. Can llms effectively leverage graph structural information: When and why. CoRR, abs/2309.16595, 2023. doi: 10.48550/ARXIV.2309. 684 16595. URL https://doi.org/10.48550/arXiv.2309.16595. 685
- 686 Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. Atlas: Few-shot learning 688 with retrieval augmented language models. Journal of Machine Learning Research, 24(251): 689 1-43, 2023. 690
- Glen Jeh and Jennifer Widom. Scaling personalized web search. In Gusztáv Hencsey, Bebo White, 691 Yih-Farn Robin Chen, László Kovács, and Steve Lawrence (eds.), Proceedings of the Twelfth In-692 ternational World Wide Web Conference, WWW 2003, Budapest, Hungary, May 20-24, 2003, pp. 693 271-279. ACM, 2003. doi: 10.1145/775152.775191. URL https://doi.org/10.1145/ 694 775152.775191.
- 696 Bowen Jin, Chulin Xie, Jiawei Zhang, Kashob Kumar Roy, Yu Zhang, Zheng Li, Ruirui Li, 697 Xianfeng Tang, Suhang Wang, Yu Meng, and Jiawei Han. Graph chain-of-thought: Aug-698 menting large language models by reasoning on graphs. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), Findings of the Association for Computational Linguistics, ACL 699 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024, pp. 163-184. Associa-700 tion for Computational Linguistics, 2024. doi: 10.18653/V1/2024.FINDINGS-ACL.11. URL 701 https://doi.org/10.18653/v1/2024.findings-acl.11.

Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick S. H. Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (eds.), Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020, pp. 6769–6781. Association for Computational Linguistics, 2020. doi: 10.18653/V1/2020.EMNLP-MAIN.550. URL https://doi.org/10.18653/v1/2020.
emnlp-main.550.

- Thomas N. Kipf and Max Welling. Variational graph auto-encoders. *CoRR*, abs/1611.07308, 2016.
 URL http://arxiv.org/abs/1611.07308.
- Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional net works. In 5th International Conference on Learning Representations, ICLR 2017, Toulon,
 France, April 24-26, 2017, Conference Track Proceedings. OpenReview.net, 2017. URL https:
 //openreview.net/forum?id=SJU4ayYgl.
- Johannes Klicpera, Aleksandar Bojchevski, and Stephan Günnemann. Predict then propagate: Graph neural networks meet personalized pagerank. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net, 2019. URL https://openreview.net/forum?id=H1gL-2A9Ym.
- Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pp. 3045–3059. Association for Computational Linguistics, 2021. doi: 10.18653/V1/2021.EMNLP-MAIN. 243. URL https://doi.org/10.18653/v1/2021.emnlp-main.243.

727 Patrick S. H. Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Na-728 man Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive NLP tasks. In 729 Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-730 Tien Lin (eds.), Advances in Neural Information Processing Systems 33: Annual Con-731 ference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 732 2020, virtual, 2020. URL https://proceedings.neurips.cc/paper/2020/hash/ 733 6b493230205f780e1bc26945df7481e5-Abstract.html. 734

- Guohao Li, Matthias Müller, Bernard Ghanem, and Vladlen Koltun. Training graph neural networks with 1000 layers. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th In- ternational Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 6437–6449. PMLR, 2021. URL
 http://proceedings.mlr.press/v139/li210.html.
- Qimai Li, Zhichao Han, and Xiao-Ming Wu. Deeper insights into graph convolutional networks for semi-supervised learning. In Sheila A. McIlraith and Kilian Q. Weinberger (eds.), Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018, pp. 3538–3545. AAAI Press, 2018. doi: 10.1609/AAAI.V32I1.11604. URL https://doi. org/10.1609/aaai.v32i1.11604.
- Rui Li, Jiwei Li, Jiawei Han, and Guoyin Wang. Similarity-based neighbor selection for graph llms. *CoRR*, abs/2402.03720, 2024a. doi: 10.48550/ARXIV.2402.03720. URL https://doi.org/ 10.48550/arXiv.2402.03720.
- Yuanchun Li, Hao Wen, Weijun Wang, Xiangyu Li, Yizhen Yuan, Guohong Liu, Jiacheng Liu, Wenxing Xu, Xiang Wang, Yi Sun, Rui Kong, Yile Wang, Hanfei Geng, Jian Luan, Xuefeng Jin, Zilong Ye, Guanjing Xiong, Fan Zhang, Xiang Li, Mengwei Xu, Zhijun Li, Peng Li, Yang Liu, Ya-Qin Zhang, and Yunxin Liu. Personal LLM agents: Insights and survey about the capability, efficiency and security. *CoRR*, abs/2401.05459, 2024b. doi: 10.48550/ARXIV.2401.05459. URL https://doi.org/10.48550/arXiv.2401.05459.

784

Sitao Luan, Chenqing Hua, Qincheng Lu, Jiaqi Zhu, Mingde Zhao, Shuyuan Zhang, Xiao-Wen Chang, and Doina Precup. Revisiting heterophily for graph neural networks. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022, 2022. URL http://papers.nips.cc/paper_files/paper/2022/hash/ 092359ce5cf60a80e882378944bf1be4-Abstract-Conference.html.

- Frxue Min, Runfa Chen, Yatao Bian, Tingyang Xu, Kangfei Zhao, Wenbing Huang, Peilin Zhao, Junzhou Huang, Sophia Ananiadou, and Yu Rong. Transformer for graphs: An overview from architecture perspective. *CoRR*, abs/2202.08455, 2022. URL https://arxiv.org/abs/2202.08455.
- Gautam Mittal, Jesse H. Engel, Curtis Hawthorne, and Ian Simon. Symbolic music generation with diffusion models. In Jin Ha Lee, Alexander Lerch, Zhiyao Duan, Juhan Nam, Preeti Rao, Peter van Kranenburg, and Ajay Srinivasamurthy (eds.), *Proceedings of the 22nd International Society for Music Information Retrieval Conference, ISMIR 2021, Online, November 7-12, 2021*, pp. 468–475, 2021. URL https://archives.ismir.net/ismir2021/paper/000058.pdf.
- Lawrence Page. The pagerank citation ranking: Bringing order to the web. Technical report, Technical Report, 1999.
- Bryan Perozzi, Bahare Fatemi, Dustin Zelle, Anton Tsitsulin, Seyed Mehran Kazemi, Rami Al-Rfou, and Jonathan Halcrow. Let your graph do the talking: Encoding structured data for llms. *CoRR*, abs/2402.05862, 2024. doi: 10.48550/ARXIV.2402.05862. URL https://doi.org/ 10.48550/arXiv.2402.05862.
- Yijian Qin, Xin Wang, Ziwei Zhang, and Wenwu Zhu. Disentangled representation learning with large language models for text-attributed graphs. *CoRR*, abs/2310.18152, 2023. doi: 10.48550/ARXIV.2310.18152. URL https://doi.org/10.48550/arXiv.2310.18152.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi
 Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to text transformer. J. Mach. Learn. Res., 21:140:1–140:67, 2020. URL http://jmlr.org/
 papers/v21/20-074.html.
- Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Gallagher, and Tina Eliassi-Rad. Collective classification in network data. *AI Mag.*, 29(3):93–106, 2008. doi: 10.1609/AIMAG. V29I3.2157. URL https://doi.org/10.1609/aimag.v29i3.2157.
- Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettle moyer, and Wen-tau Yih. REPLUG: retrieval-augmented black-box language models. CoRR,
 abs/2301.12652, 2023. doi: 10.48550/ARXIV.2301.12652. URL https://doi.org/10.
 48550/arXiv.2301.12652.
- Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. Advances in neural information processing systems, 30, 2017.
- Yanchao Tan, Hang Lv, Xinyi Huang, Jiawei Zhang, Shiping Wang, and Carl Yang. Musegraph: Graph-oriented instruction tuning of large language models for generic graph mining. *CoRR*, abs/2403.04780, 2024. doi: 10.48550/ARXIV.2403.04780. URL https://doi.org/10. 48550/arXiv.2403.04780.
- Jiabin Tang, Yuhao Yang, Wei Wei, Lei Shi, Lixin Su, Suqi Cheng, Dawei Yin, and Chao Huang.
 Graphgpt: Graph instruction tuning for large language models. In Grace Hui Yang, Hongning
 Wang, Sam Han, Claudia Hauff, Guido Zuccon, and Yi Zhang (eds.), Proceedings of the 47th
 International ACM SIGIR Conference on Research and Development in Information Retrieval,
 SIGIR 2024, Washington DC, USA, July 14-18, 2024, pp. 491–500. ACM, 2024. doi: 10.1145/
 3626772.3657775. URL https://doi.org/10.1145/3626772.3657775.

810	Heng Wang, Shangbin Feng, Tianxing He, Zhaoxuan Tan, Xiaochuang Han, and Yulia Tsvetkov.
811	Can language models solve graph problems in natural language? In Alice Oh, Tris-
812	tan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), Ad-
813	vances in Neural Information Processing Systems 36: Annual Conference on Neural Infor-
814	mation Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16,
815	2023, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/
816	622afc4edf2824a1b6aaf5afe153fa93-Abstract-Conference.html.

- Yanbang Wang, Hejie Cui, and Jon M. Kleinberg. Microstructures and accuracy of graph recall by large language models. *CoRR*, abs/2402.11821, 2024. doi: 10.48550/ARXIV.2402.11821. URL https://doi.org/10.48550/arXiv.2402.11821.
- Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. Finetuned language models are zero-shot learners. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022.* OpenReview.net, 2022. URL https://openreview.net/forum?id=gEZrGCozdqR.
- Michihiro Yasunaga, Armen Aghajanyan, Weijia Shi, Richard James, Jure Leskovec, Percy Liang, Mike Lewis, Luke Zettlemoyer, and Wen-Tau Yih. Retrieval-augmented multimodal language modeling. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *International Conference on Machine Learning, ICML 2023,* 23-29 July 2023, Honolulu, Hawaii, USA, volume 202 of Proceedings of Machine Learning Research, pp. 39755–39769. PMLR, 2023. URL https://proceedings.mlr.press/ v202/yasunaga23a.html.
- Ruosong Ye, Caiqi Zhang, Runhui Wang, Shuyuan Xu, and Yongfeng Zhang. Natural language
 is all a graph needs. *CoRR*, abs/2308.07134, 2023. doi: 10.48550/ARXIV.2308.07134. URL
 https://doi.org/10.48550/arXiv.2308.07134.
- Wenhao Yu, Dan Iter, Shuohang Wang, Yichong Xu, Mingxuan Ju, Soumya Sanyal, Chenguang Zhu, Michael Zeng, and Meng Jiang. Generate rather than retrieve: Large language models are strong context generators. In *The Eleventh International Conference on Learning Representations, ICLR* 2023, Kigali, Rwanda, May 1-5, 2023. OpenReview.net, 2023. URL https://openreview. net/forum?id=fB0hRu9GZUS.
- Delvin Ce Zhang, Menglin Yang, Rex Ying, and Hady W. Lauw. Text-attributed graph representation learning: Methods, applications, and challenges. In Tat-Seng Chua, Chong-Wah Ngo, Roy Ka-Wei Lee, Ravi Kumar, and Hady W. Lauw (eds.), *Companion Proceedings of the ACM on Web Conference 2024, WWW 2024, Singapore, Singapore, May 13-17, 2024*, pp. 1298–1301. ACM, 2024a. doi: 10.1145/3589335.3641255. URL https://doi.org/10.1145/3589335. 3641255.
- Jiawei Zhang. Graph-toolformer: To empower llms with graph reasoning ability via prompt augmented by chatgpt. *arXiv preprint arXiv:2304.11116*, 2023.
- Zeyang Zhang, Xin Wang, Ziwei Zhang, Haoyang Li, Yijian Qin, and Wenwu Zhu. Llm4dyg: Can large language models solve spatial-temporal problems on dynamic graphs? In Ricardo Baeza-Yates and Francesco Bonchi (eds.), *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD 2024, Barcelona, Spain, August 25-29, 2024*, pp. 4350–4361. ACM, 2024b. doi: 10.1145/3637528.3671709. URL https://doi.org/10. 1145/3637528.3671709.
- Jianan Zhao, Meng Qu, Chaozhuo Li, Hao Yan, Qian Liu, Rui Li, Xing Xie, and Jian Tang. Learning on large-scale text-attributed graphs via variational inference. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023.* OpenReview.net, 2023a. URL https://openreview.net/forum?id=q0nmYciuuZN.

Jianan Zhao, Le Zhuo, Yikang Shen, Meng Qu, Kai Liu, Michael M. Bronstein, Zhaocheng Zhu, and Jian Tang. Graphtext: Graph reasoning in text space. *CoRR*, abs/2310.01089, 2023b. doi: 10.
 48550/ARXIV.2310.01089. URL https://doi.org/10.48550/arXiv.2310.01089.

A Architecture and Training of the Graph Neural Network ψ

In our setting, semi-supervised node classification problem, $\mathcal{V}, \mathcal{E}, \mathcal{T}, \mathcal{Y}^{\text{train}}$ are available during training. Since Graph Neural Networks (GNNs) are not inherently capable of processing textual features, we employ a pretrained text encoder to generate *d*-dimensional dense embeddings for each node

$$\operatorname{Encoder}_{\psi_1}(t_i) = \mathbf{h}_i^{(0)} \in \mathbb{R}^d, \forall i \in 1, \dots, n$$
(15)

Subsequently, we apply a standard graph neural network. For this study, we adopt Graph-SAGE (Hamilton et al., 2017) whose iterative architecture is

$$\mathbf{h}_{i}^{(l)} \leftarrow \sigma^{(l)} \Big(\operatorname{MEAN}\Big(\{ \mathbf{h}_{i}^{(l-1)} \} \cup \{ \mathbf{h}_{j}^{(l-1)} : (v_{i}, v_{j}) \in \mathcal{E} \} \Big) \cdot \mathbf{W}^{(l)} \Big)$$
(16)

where $\sigma^{(l)}$ is the activation function and $\mathbf{W}^{(l)}$ is the learnable parameter of each layer. For an *L*-layed network, in the last layer, $\sigma^{(L)}$ is selected as Softmax and $\mathbf{W}^{(L)} \in \mathbb{R}^{d \times c}$ resulting in $\mathbf{h}_i^{(L)} \in \mathbb{R}^c$ as the prediction vector. The typical loss used for training the GNN is negative loglikelihood $\mathcal{L}_{\text{NLL}}(\mathbf{h}_i^{(L)}, y_i)$ for all the nodes in the training set $\mathcal{Y}^{\text{train}}$. The complete set of trainable parameters is denoted as $\psi = \{\psi_1\} \cup \{\mathbf{W}^{(l)}\}_{l=1}^L$.

B INTERPRETATION OF THE DISTRIBUTION MATCHING LOSS

We recap the objective function. For notation brevity, we use t_i to denote the input target node t^{target} :

 $\mathrm{KL}(\tilde{p}_{\mathrm{LM}}(\cdot|t_i, y_i) \| p_{\phi}(\cdot|t_i)) \tag{17}$

where the stop gradient operator is removed if we only compute gradient with respect to ϕ and

$$p_{\phi}(t_j|t_i) = \frac{\exp(s_{\phi}(t_i, t_j))}{\sum_{t_k \in \mathcal{D}} \exp(s_{\phi}(t_i, t_k))}$$
(18)

and

$$\tilde{p}_{\mathrm{LM}}(t_j|t_i, y_i) = \frac{\exp(p_{\mathrm{LM}}(y_i|t_i, t_j))}{\sum_{k \in \mathcal{N}_i} \exp(p_{\mathrm{LM}}(y_i|t_i, t_k))}$$
(19)

For notation brevity, we replace $\sum_{t_{k} \in D}$ with \sum_{z} if there is no ambiguity. Then

$$\min_{\phi} \operatorname{KL}(\tilde{p}_{\mathsf{LM}}(\cdot|t_i, y_i) \| p_{\phi}(\cdot|t_i)) \Leftrightarrow \min_{\phi} - \sum_{z} \tilde{p}_{\mathsf{LM}}(z|t_i, y_i) \log[p_{\phi}(z|t_i)]$$
(20)

$$= -\sum_{z} \tilde{p}_{\text{LM}}(z|t_i, y_i) \log\left[\frac{e^{s_{\phi}(z, t_i)}}{\sum_{z'} e^{s_{\phi}(z', t_i)}}\right]$$
(21)

$$=\sum_{z} \tilde{p}_{\text{LM}}(z|t_i, y_i) \log\left[\sum_{z'} e^{s_{\phi}(z', t_i)}\right] - \sum_{z} \tilde{p}_{\text{LM}}(z|t_i, y_i) s_{\phi}(z, t_i)$$
(22)

$$= \log \left[\sum_{z} e^{s_{\phi}(z,t_i)} \right] - \sum_{z} \tilde{p}_{\text{LM}}(z|t_i, y_i) s_{\phi}(z,t_i)$$
(23)

909 Hence,

$$\nabla \mathrm{KL} = \frac{\sum_{z} e^{s_{\phi}(z,t_i)} \nabla s_{\phi}(z,t_i)}{\sum_{z'} e^{s_{\phi}(z',t_i)}} - \sum_{z} \tilde{p}_{\mathrm{LM}}(z|t_i,y_i) \nabla s_{\phi}(z,t_i)$$
(24)

$$=\sum \left[p_{\phi}(z|t_i) - \tilde{p}_{\text{LM}}(z|t_i, y_i) \right] \nabla s_{\phi}(z, t_i)$$
(25)

917
$$= \sum_{z} \left[1 - \frac{p_{\text{LM}}(z|t_i, y_i)}{p_{\phi}(z|t_i)} \right] p_{\phi}(z|t_i) \nabla s_{\phi}(z, t_i)$$
(26)

Temp	plate name		Text							
Citation (for Cora, Pubmed, ogbn-arxiv) Citation Title Last (for Cora, Pubmed, ogbn-arxiv)			Please classify the following paper into {pruned label candidates} based on the provided information\nTitle: {target node's title}\nContent: {target node's abstract}\nRelated papers: {retrieved nodes' titles}Please classify the following paper into {pruned label candidates} based on the provided information\nContent: {target node's abstract}\nRelated papers: {retrieved nodes' titles}Please classify the following paper into {pruned label candidates} based on the provided information\nContent: {target node's abstract}\nRelated papers: {retrieved nodes' titles}\nTitle: {target node's title}Please classify the following Amazon product into {pruned label candidates} based on the provided information\nProduct name: {target node's title}\nDescription: {target node's description}\nRelated products: {retrieved nodes' titles}Please classify the following Amazon product into {pruned label candidates} based on the provided information\nProduct into {pruned label candidates} based on the provided information\nProduct into {pruned label candidates} based on the provided information\nProduct into {pruned label candidates} based on the provided information\nProduct into {pruned label candidates} based on the provided information\nProduct into {pruned label candidates} based on the provided information\nProduct into {pruned label candidates} based on the provided information\nProduct into {pruned label candidates} based on the provided information\nProduct into {pruned label candidates} based on the provided information\nProduct into {pruned label candidates} based on the provided information\nProduct into {pruned label candidates} based on the provided information\nProduct into {pruned label candidates} based on the provided information\nProduct into {products: {retrieved nodes' titles}\nProduct name: {target node's title}Table 5: Dataset statistics.							
							Amazon (for ogbn-products)			
Amazon Title Last (for ogbn-products)										
							_	Name	# nodes	# edges
_	Cora	2,708					10,556	7	Random 60/20/20%	Accuracy
	Pubmed	19717	88 648	3	Random 60/20/20%	Accuracy				

Table 4: Templates used for all datasets.

After changing the notation back from \sum_{z} to $\sum_{t_k \in \mathcal{D}}$, we have

169,343

54,025

1,166,243

198,663

$$\nabla \mathrm{KL} = \sum_{t_k \in \mathcal{D}} \left[1 - \frac{\tilde{p}_{\mathrm{LM}}(t_j | t_i, y_i)}{p_{\phi}(t_j | t_i)} \right] p_{\phi}(t_j | t_i) \nabla s_{\phi}(t_j, t_i)$$
(27)

Given split

Given split

Accuracy

Accuracy

whose rationale is that if the LM's feedback greatly prefers the neighbor v_i (and its associated text t_j), larger than its probability to be retrieved by the retriever (i.e., $\frac{\tilde{p}_{LM}(t_j|t_i,y_i)}{p_{\phi}(t_j|t_i)} > 1$), then the similarity score between t_i and t_j will increase, i.e., improve the probability of t_j to be retrieved.

C TEMPLATES

ogbn-arxiv

ogbn-products

Table 4 presents templates used in this paper. We design the "Citation" template for the Cora, Pubmed, and ogbn-arxiv datasets and the "Amazon" template for the ogbn-products dataset.

Drawing inspiration from the findings of He et al. (2024), who demonstrated the efficacy of positioning the title after the main content for certain datasets, we have also introduced two additional template variations: "Citation Title Last" and "Amazon Title Last."

D **DATASET STATISTICS**

We present the detailed statistics of datasets used in this paper in Table 5.

972 E ADDITIONAL EXPERIMENTS

974 E.1 ADDITIONAL EFFICIENCY STUDY

FLOPs Here we compare the floating point operations (FLOPs) of our proposed AUGGLM and our main baseline, InstructGLM Ye et al. (2023), a text-output node classifier.

The computation of InstructGLM includes (1) computing node encodings, which is precomputed, and (2) training and inference of the downstream LM. In contrast, the computation of our AugGLM includes (1) computing PPR neighbors for every node, which is also precomputed, (2) training and inference of the semantic retriever, and (3) training and inference of the downstream LM. Hence, the extra on-the-fly computation overhead comes from the semantic retriever, all-MiniLM-L6-v2, in our experiments. We report the FLOPs of the retriever and different LM backbones in Table 6.

Table 6:	FLOPs	comparison	between	different	modules.
rable 0.	LOIS	companison	UCC WCCII	uniterent	mouules.

987		FLOPs
988	Retriever	2.3G
989	FLAN-T5 (small)	71.7G
990	FLAN-T5 (base)	257.2G
991	FLAN-T5 (large)	845.4G
000		

The results show that the retriever only adds a tiny amount of FLOPs compared to the backbone LMs. In other words, the FLOPs of our framework are very close to those of InstructGLM if we adopt the same downstream LM. More concretely, if both our model and InstructGLM select T5-large as the backbone, the FLOPs of InstructGLM would be 845.4G, and our framework would be 847.7G.

Memory usage. Memory usage is linear concerning batch size. We report the memory usage with different backbone LMs in Table 7, where we set the batch size to 1.

Table 7: Memory usage with different LMs.

	GPU memory	
AugGLM (small) AugGLM (base) AugGLM (large)	3098M 6572M 20308M	

1010 It is reasonable that more powerful backbone LMs require more GPU memory.

Convergence analysis. We train AUGGLM with different backbones: FLAN-T5small/base/large on the Cora dataset and plot their loss curve regarding updating steps in
Figure 4. It shows that our proposed AUGGLM is easy to train and converges smoothly when
equipped with various backbones LMs of different scales.

Running time. We recorded the running time (both forward and backpropagation) of the semantic retriever and the backbone LMs in Table 8. The dataset we tested was Cora, and the batch size was 1. Note that this wall clock running time is related to the batch size, dataset, and specific hardware. In this experiment, the running time is tested on an NVIDIA A100-SXM4-80GB.

1021 Overall, we can conclude that the semantic retriever only adds very limited on-the-fly computation overhead compared to the downstream LM, showing the efficiency of our proposed framework.

1024 E.2 Additional Parameter Study

In this section, we study the model's performance with various hyperparameters.



Figure 4: Convergence curve of AUGGLM.

Table 8: Wall-clock running time comparison between different modules.

	Forward (ms)	Backprop (ms)
Retriever	14.7	6.1
FLAN-T5 (small)	90.0	32.0
FLAN-T5 (base)	104.4	66.6
FLAN-T5 (large)	277.2	197.0

Selection of the backbone GNN. Specifically, we study the performance of AUGGLM equipped with different GNNs. we compared the performance of AugGLM equipped with GraphSAGE (used in the reported results) with the counterpart equipped with GCN Kipf & Welling (2017). The comparison is in Table 9.

Table 9: Performance comparison of AUGGLM equipped with different GNNs.

	Cora	Pubmed	ogbn-arxiv	ogbn-products
GraphSAGE	91.14	94.80	75.39	81.73
GCN	90.98	94.85	75.21	81.82

We observed that the performance is nearly identical between GCN and GraphSAGE. This can be attributed to two factors: (1) the classification performances of GCN and GraphSAGE are similar, and (2) the GNN is used to generate prototypes and prune candidate labels, which does not require a highly powerful GNN for accurate classification.

Number of PPR retrieved nodes. Next, we preliminarily examined the relationship between the model performance and the number of nodes retrieved. In this auxiliary experiment, we fixed the number of nodes retrieved by semantic retrieval at 5 and varied the number of nodes retrieved by PPR retrieval. The results are reported in Table 10

Table 10: Performance comparison of AUGGLM with different PPR retrieved neighbors.

PPR neighbors	1	3	5	7	9	10	15	20	25
ogbn-arxiv	75.18	75.76	75.39	75.19	76.05	76.45	75.99	74.81	74.48

1076 Interestingly, we found that the model's performance remains relatively stable when the number of
1077 PPR nodes is less than 15. However, the performance degrades when too many nodes are retrieved
1078 (more than 15). A possible explanation is that when the number of PPR nodes becomes too large,
1079 every target node's retrieved nodes become similar (e.g., some hub nodes are retrieved by most nodes), reducing the discriminativeness of each target node. This phenomenon is reminiscent of the

"oversmoothing" problem Li et al. (2018) in GNNs, where a GNN with too many layers produces indistinguishable latent representations for all nodes.

Other topological retrieval options. In this auxiliary experiment, we use the link predictor to retrieve relevant neighbors. Specifically, we trained a graph autoencoder (GA) Kipf & Welling (2016), a basic graph neural network-based link predictor, on the given graph. Then, we retrieved the top-5 most confident neighbors from the reconstructed graph to replace those obtained through PPR retrieval* The results are presented in Table 11, where Flan-T5-small is used as the backbone. For better reference, we also provide a version where PPR retrieval is replaced with retrieving from 1-hop neighbors.

Table 11: Performance comparison of AUGGLM with different topological retrieval techniques.

	Cora	Pubmed	ogbn-arxiv	ogbn-products
1-hop neighbors	90.59	94.33	73.97	79.53
GA	90.83	94.42	74.01	79.85
PPR neighbors	91.14	94.80	75.39	81.73

We observe that both 1-hop neighbor retrieval and GA perform worse than their PPR counterparts.
A possible reason is that both 1-hop neighbor retrieval and GA are local retrieval methods, whereas
PPR can effectively capture the global structure. Additionally, we note that GA is trained using
a reconstruction loss, which means it tends to assign high confidence to existing edges. In other
words, the neighbors retrieved by GA would be similar to those obtained through 1-hop neighbor
retrieval, except for some low-degree nodes.

1105 F SELECTED HYPERPARAMETERS

1107 We report the hyperparameter used for every dataset in Table 12. More detailed hyperparameters1108 will be released with the code upon publication.

	Cora	Pubmed	ogbn-arxiv	ogbn-products
# PPR neighbors	5	2	5	5
# Semantic neighbors	5	2	5	5
Template	Citation	Citation	Citation title last	Amazon
# Candidate labels	3	2	3	3
LM learning rate	1e-4	1e-4	1e-4	1e-4
Retriever learning rate	1e-5	1e-5	1e-5	1e-5
Weight decay	0	0	0	0

Table 12: Hyperparameters selected of AUGGLM.