# LINK PREDICTION ON TEXT ATTRIBUTED GRAPHS: A NEW BENCHMARK AND EFFICIENT LM-NESTED GRAPH CONVOLUTION NETWORK DESIGN

Anonymous authors

Paper under double-blind review

### ABSTRACT

Textual and topological information is significant for link prediction (LP) in textattributed graphs (TAGs). Recent link prediction methods have focused on improving the performance of capturing structural features by Graph Convolutional Networks (GCNs). The importance of enhancing text encodings, powered by the advanced Pre-trained Language Models (PLM) has been underestimated. In this work, we analyse and emphasise the importance of PLMs and propose a novel PLM-nested GCN design. We developed an extensive benchmark to compare current competitive link prediction methods and PLM-based methods in a unified experimental setting and systematically investigate the representation power of the text encoders in the link prediction task. Based on our investigation, we introduce LMGJOINT — a memory-efficient fine-tuning method. The key design features include: residual connection of textual proximity, a combination of structural and textual embeddings and a cache embedding training strategy. Our empirical analysis shows that these design elements improve MRR by up to 19.75% over previous state-of-the-art methods and achieve competitive performance across a wide range of models and datasets.

027 028 029

006

007

008 009 010

011

013

014

015

016

017

018

019

021

024

025

026

## 1 INTRODUCTION

Link Prediction (LP) aims to predict the likelihood of a connection between two nodes in a graph,
encompassing various real-world applications, including protein-protein interaction prediction (Szklarczyk et al., 2018), recommendation systems (Huang et al., 2005) and knowledge graph completion
(Hu et al., 2020b). While early LP relied on handcrafted graph heuristics (Adamic & Adar, 2003),
more advanced approaches follow a two-stage framework: (1) an encoder maps graph information
into node embeddings and (2) then a decoder assesses pairwise embedding similarity to predict
connection likelihood.

Among encoder designs, Graph Convolutional Networks (GCNs) are the dominant paradigm, de pending on both node and structural features. In previous benchmarks such as Cora (McCallum
 et al., 2000) and PubMed (Sen et al., 2008b), node features have often relied on shallow text em beddings such as Word2Vec (Mikolov et al., 2013). However, these embeddings struggle to capture context-aware information and complex semantic relationships, which are crucial for link prediction.

Despite their limitations, these shallow embeddings are widely used in standard benchmarks (Hu et al.) [2021], which has led to several issues: (1) they are often practically irreproducible, making it difficult to replicate them on new datasets; (2) the reliance on a specific text embedding has resulted in architecture over-optimization in current algorithm development; and (3) the decoupling of the text embedding process and the GCN design hinders seamless end-to-end training, thereby reducing the overall effectiveness of the approach.

Text-attributed graphs (TAGs) help overcome these limitations by offering rich semantic content.
They enable the characterization of individual node properties using powerful Pretrained Language
Models (PLMs). Additionally, TAGs allow for the seamless integration of learnable text embeddings
with GCN-based structural aggregation. However, existing works on TAGs primarily focus on node
classification (Duan et al., 2024; He et al., 2023; Yang et al., 2021; Zhu et al., 2024b) or suffer
from limited comparisons due to a lack of strong baselines (Wang et al., 2023; Yun et al., 2021;



Figure 1: The overview of LMGJOINT. The framework consists of three main components: (a) structure embedding  $\mathbf{H}_{C}$  (soft/hard common neighbors) from the adjacency matrix A. (b)  $\mathbf{S}^{T}$ semantic embeddings proximity based on sentence embedding derived from PLM. (c) Aggregated embeddings  $\mathbf{H}^{K}$  which incorporate both  $\mathbf{X}^{T}$  and  $\mathbf{X}$ . (d) The final step concatenates (a,b,c) through a MLP to generate link prediction.

Chamberlain et al., 2023; Zhang et al., 2020; Zhang & Chen, 2018). Motivated by this limitation, we make the following contributions:

- 1.Data contribution We collect and introduce ten graphs including Small PaperwithCode (Saier 076 et al.] 2023), Cora (McCallum et al.] 2000), Arxiv\_2023 (He et al., 2023), PubMed (Sen et al. 077 2008b), Medium PaperwithCode (Saier et al., 2023), Photo Shchur et al. (2018), History (Li et al., 2024), Ogbn-arxiv (Hu et al., 2021), Citationv8 (Wu et al., 2021) and Ogbn-papers100M (Hu et al. 2021). We provide rich statistics compactly describing their density, hierarchy, locality and generalized node homophily. These datasets and statistics serve as a foundation for exploring these new hypotheses driving the link prediction community moving forward.
- 082 2.Extensive Empirical Benchmark Using the proposed datasets, we have provided a thorough 083 benchmark, offering a fair comparison of ten GCN-based link prediction approaches alongside 084 seven traditional path-based methods. These selections broadly represent the current LP algorithm 085 space, including state-of-the-art methods. Additionally, we expand the PLM-based baselines by adapting analogous architectures from node classification. This includes both cascade and nested 087 architectures, as discussed in Section 5. Our benchmark is available at TAG4LP.
- 3.LMGJOINT, a powerful nested framework We introduce Language Model Graph Joint Design (LMGJOINT), a parameter- and memory-efficient method. We identify three key design features, 090 including (D1) residual connection of textual proximity, (D2) combination of structural and textual 091 embeddings and (D3) cache embedding training strategy. We provide a theoretical justification for these design principles in Section 4. Our integration results in a fine-tuned architecture 092 that preserves the GCN's strength in structural feature extraction while leveraging the PLM's ability to capture complex semantic relationships. Experimental comparisons with state-of-the-art 094 approaches demonstrate that LMGJOINT achieves up to a 19.75% improvement in MRR. Moreover, 095 experiments across seven proposed datasets, scaling up to  $10^7$  nodes, validate the effectiveness 096 and scalability of our approach, consistently outperforming competitive baselines. Furthermore, 097 LMGJOINT is not limited to specific GCNs. It can be easily combined with any graph-based 098 model and PLM-based text encoder without requiring any changes to the latter or affecting its 099 computational complexity.
- 100 101 102

103

054

056

060

061

062

063

064

065 066

067

068

069

070

071 072 073

074

075

079

081

#### **RELATED WORK** 2

104 LM-based approaches for TAGs. Shallow embedding: In the context of TAGs, previous preprocess-105 ing methods often involved transforming text attributes into bag-of-words (Harris, 1954). Though widely adopted in the graph community, it has limited capacity to capture complex semantic relation-106 ships or fully utilize the richness of text attributes provided by modern pre-trained language models 107 (PLMs). PLM-based method: To address these limitations, recent approaches leverage fine-tuning of

2

108 109

110

111

112

113

114

115

116

117 118

119

120

124 125



Figure 2: The figure compares TAG4LP (yellow) with previous ogbl datasets (blue) (Hu et al., 2021). Circle sizes indicate generalized edge homophily (Eq 44), with non-featured graphs (gray) set to 0.5. For more details, see Appendix G

pre-trained LMs to generate node embeddings tailored to the domain and context of TAGs. It can 126 be classified into two main frameworks: (1) Cascade framework: Text embedding from PLM and 127 graph aggregation are performed sequentially without interaction. Examples include SimTeG (Duan 128 et al. 2024), GAINT (Chien et al. 2021) and TAPE (He et al. 2023). (2) Nested framework: PLMs 129 and GCNs are optimized jointly, enabling iterative or integrated learning. For instance, Graphformer 130 integrates text encoding and graph aggregation in an iterative workflow (Yang et al., 2021), while 131 Engine incorporates caching and a dynamic early-exit mechanism to enhance performance and reduce 132 training costs in Llama 3 (Zhu et al., 2024b). A detailed comparison of our benchmark with related 133 work is provided in Table 1. (3) Instruction Learning: GraphGPT (Tang et al., 2023) integrates 134 LLMs with graph structural knowledge through instruction tuning. Furthermore, LinkGPT (He et al., 135 2024) proposes a two-stage fine-tune method, achieving state-of-the-art performance in zero-shot and 136 few-shot settings.

137 **Related Benchmarks** Our proposed method bears methodological resemblance to (Zhang et al., 138 2024a). It benchmarks a co-training method on 22 graphs in both link prediction and node classi-139 fication tasks. However, the proposed approach fails to bring performance gain for link prediction. 140 Besides, Mao et al. (2023) critically examines the fundamental incompatibility between node features 141 and structural similarity, which grounds the analysis from a data science perspective. Li et al. (2023) 142 proposes a benchmark of all existing GCN4LP methods under consistent data splits and training settings. Their findings reveal that advancements in GCN4LP are primarily due to improved capture 143 of pairwise structural features. Similarly, Wu et al. (2021) introduces a new TAG benchmark, mainly 144 focusing on node classification. It proposes a co-training paradigm by simply concatenating GNNs 145 and LLM/PLMs several cascade-architectures without a task-specialised design. Recent work starts 146 considering including edge textual features into TAG and conducting various experiments on cascade 147 GCN-LLM models (Li et al.) 2024). In summary, while current methods for node classification have 148 provided foundational insights into pretraining tasks and algorithm design, there is no counterpart 149 and established SOTA for link classification tasks. 150

**Related Datasets** The widely used datasets for link prediction (LP) were introduced by OGB (Hu 151 et al., 2021), but their rich textual attributes have been largely underexplored. Recently, TAPE (He 152 2023) and TAG\_Benchmark (Yan et al., 2023) introduced several text-attributed graphs (TAGs), et al. 153 such as Cora (McCallum et al.) 2000), PubMed (Sen et al.) 2008b), and Arxiv\_2023 (He et al.) 154 2023). The Engine further expanded these datasets into seven graphs. Similarly, Chen et al. (2024b) 155 conducted preliminary studies on three datasets for LP. It observes that without task-specific design 156 simply combining LLM and GCN in a nested architecture fails to achieve performance improvements. 157 Recently, edge-level textual features have garnered attention, with Li et al. (2024) introducing a 158 benchmark of 9 graphs. However, this benchmark is limited to a cascade GCN-LLM design. Building 159 on insights from these works, our benchmark focuses initially on homophilic networks (Hu et al., 2021) and later generalizes to other domains and non-attributed graphs. Key distinctions from existing 160 benchmarks are highlighted in Figure 2. In summary, our benchmark is the most comprehensive, 161 featuring the widest variety of algorithms and the largest number of datasets evaluated.

#### 162 **DATASET CONTRIBUTION** 3

163 164

Data factors and Current Limitation: The further 165 development of the Link Prediction algorithm is hin-166 dered by the efficient hypothesis. The limitations 167 of applying GNNs for node classification on het-168 erophily graphs are well understood. In comparison, 169 prior works on GNN4LP are mostly based on hand-170 crafted structure features (Zhang et al., 2020; Zhang & Chen, 2018; Wang et al., 2023; Yun et al., 2021). 171 Despite the practical improvement, our understand-172 ing of the dominant data factor within GNN4LP 173 remains incomplete. We identify three critical data 174 factors: 1) Feature homophily refers to the impact 175 of similar features, a recent study indicates discrep-176 ancies between feature proximity and structure Zhu 177 et al. (2024a) leads to performance decay for link

Table 1: Comparison with existing methods.  $\checkmark$ : public benchmark,  $\checkmark$ : benchmarked model, NC: node classification, LP: link prediction, Num: number of the evaluated dataset

Tack	Works	Cas	cade	N	Num		
Task	WOIKS	MLP	GCN	MLP	GCN	Itum	
NC	Graphformer				1	3	
NC/LP	GAINT	1	1			2	
NC/LP	SimTeG	1	1			3	
NC	TAPE		1			5	
NC	LEADING		1			3	
NC	ENGINE				1	7	
NC	TEG-DB	1	1			9	
LP	Ours	1	1	1	1	9	

178 prediction. We quantify this data factor by generalized edge homophily defined in Appendix  $\overline{G}(2)$ 179 Structure hierarchy describes the hierarchical structure that widely exists in the citation network. 180 When embedding such a graph in Euclidean space, GCN-based embedding incurs a large distortion 181 compared to in hyperbolic space Liu et al.) Chami et al. (2019). We quantify such hierarchy using 182  $\alpha$  in degree distribution; 3) Pairwise local structure: This hypothesis originates from the intrinsic permutation invariance of GCN. It results in the limited expressivity to distinguish automorphic nodes 183 Chamberlain et al. (2023). To analyze and study their impact on link prediction from a data-centric 184 perspective, we suggest clustering and transitivity to measure such local distance features. To sum 185 up we propose 12 graph statistics, covering three categories, as shown in Table 5. These statistics compactly and thoroughly quantify the above-mentioned three data factors. Our proposed dataset and 187 statistics provide valuable resources to advance research in the TAG and GRL communities. Dataset 188 statistics can be found in Appendix G.3 189

**Proposed dataset** To address the above limitations, we introduce a novel TAG dataset comprising 190 eight graphs from prior literature and two generated graphs. This collection offers several distinct 191 advantages: (1) Expanded Scale: Our collection builds on widely adopted benchmarks in the graph 192 research community, such as Cora, PubMed and Arxiv\_2023 to ensure consistency and comparability 193 with existing studies. Additionally, we introduce two datasets derived from the PaperswithCode 194 API (Saier et al., 2023). It provides a continuous spectrum of node sizes ranging from  $10^2$  to  $10^9$ . 195 To further enhance scalability and enable the study of large-scale settings, we include Citationv8 196 (Yan et al., 2023) and ogbn-paper100M (Hu et al., 2021). (2) Enriched Textual Information: Except 197 traditional node features derived from shallow embedding (Mikolov et al.), [2013; [Harris, [1954]), our dataset also retains the original textual content associated with nodes. This enriched textual information enables more algorithmic flexibility to provide more advanced text encoding using 199 PLMs.(3) Extensive and Comprehensive Statistics: Previous datasets have typically reported only the 200 number of nodes and edges, offering limited insights into the underlying structural complexities. We 201 offer a richer set of statistics in Appendix G relates to current hypothesis, including density, hierarchy, 202 locality, and feature homophily (Zhu et al., 2024a). We illustrate the difference between proposed 203 dataset with OGB (Hu et al., 2021) in terms of scale and feature homophily in Figure 2 More details 204 about data statistics can be found in Appendix G.3. 205

206 207

208

#### 4 LMGJOINT: A NEW EFFICIENT MODEL

209 We begin by describing three basic components (C1, C2, C3) that guided our design and then highlight 210 three efficient designs (D1, D2, D3) that can help improve the performance and reduce memory 211 requirements. 212

C1: Soft Common Neighbor: On the other extreme, we utilize Common Neighbor, a first-order 213 structure feature that solely utilizes graph topology. 214

215

$$\mathbf{H}_{ij}^{\mathsf{C}} = \mathbb{I}\left(\mathbf{A}\mathbf{A}_{(i,j)} > d\right) \tag{1}$$

where  $A \in \{0, 1\}^{n \times n}$  is the binary adjacency matrix, I is indication function, d is a hard threshold to remain stronger connections. Common neighbors can also be leveraged to directly detect the likelihood of a connection between nodes, i.e. Hard Common Neighbors (Newman, 2001).

**C2:** Semantic Feature Proximity: In a homophilic setting, connected nodes exhibit high textual proximity. Thus, a straightforward approach is to disregard graph structure and train a multilayer perceptron (MLP) solely on the text encodings. Let  $\mathbf{T}, \mathbf{X}^T$  represent the raw text and embedded text features from PLM. The semantic proximity between node pair (i, j) is defined as:

$$\mathbf{X}^T = \mathsf{PLM}(\mathbf{T}) \tag{2}$$

$$\mathbf{S}_{ii}^T = \mathbf{W}(\mathbf{X}_i^T \odot \mathbf{X}_i^T) \tag{3}$$

Here W is a learned weight matrix. The operator ⊙ denotes the Hadamard product. PLM is the
pre-trained embedding model that maps raw text to a numeric vector. We benchmark three different
sentence embedding methods including e5-large-v2 (Wang et al., 2022), Sentence-Transformers
MiniLM-L6-v2 (Reimers & Gurevych, 2019a) and MPNet(Song et al., 2020).

C3: Aggregated Semantic Feature with Self-loop: The aggregated features are propagated with a self-loop to capture the information of k-step neighbors. This is useful to capture the information of similar neighbors when the structure exhibits homophily e.g. in a citation graph (Lee et al., 2024).

$$\mathbf{H}^{k} = f\left(\tilde{\mathbf{A}}_{\text{sym}}\mathbf{H}^{k-1}\mathbf{W}\right) \tag{4}$$

 $\mathbf{H}^{0} = \mathbf{X}$  and successively optimized by  $\mathbf{X}^{T}$  from PLM by cache embedding strategy introduced in Section 4.1 We symmetrically normalize the adjacency matrix  $\tilde{\mathbf{A}}_{sym} =$  $(\mathbf{D} + \mathbf{I})^{-\frac{1}{2}}(\mathbf{A} + \mathbf{I})(\mathbf{D} + \mathbf{I})^{-\frac{1}{2}}$ ,  $\mathbf{I}$ ,  $\mathbf{D}$  are the identity and diagonal degree matrix (Kipf & Welling 2016). We can collapse the repeated multiplication with the normalized adjacency matrix  $\tilde{\mathbf{A}}_{sym}$  into a single matrix to the K-th power,  $\tilde{\mathbf{A}}_{sym}^{K}$ . Then we have the aggregated feature as  $\mathbf{H}^{K} = f(\tilde{\mathbf{A}}_{sym}^{K}\mathbf{X}\mathbf{W})$ .

4.1 EFFICIENT NESTED ARCHITECTURE

224

225 226

231

232

237

238

239

240

241 242 243

244

248

Ours: LMGJOINT We combine embeddings from these three simple components through simple linear transformations and component-wise non-linearities Lee et al. (2024); Wang et al. (2023);
 Chamberlain et al. (2023).

$$\mathbf{Y} = \mathsf{MLP}\left([\mathbf{H}^{K}; \beta \mathbf{H}_{\mathsf{C}}; \mathbf{S}^{T}]\right) \tag{5}$$

**Y** is our model's output predictions.  $\beta$  is the weight for the structure feature. In Fig. 1 we visualize LMGJOINT. We also give its pseudocode in Algorithm 1. We then

251 D1: Jumping Connection of Textual Similarity GNNs primarily depend on node-level aggregation 252 via summing the weighted neighboring features iteratively. Such a local smoothing mechanism helps 253 generate more representative embedding when the homophily assumption holds Luan et al. (2023). 254 However, it also smooths the rich semantic embedding from textual nuance during the smoothing 255 process. We address such issues by combining pairwise semantic proximity without training at the last layer, i.e., the semantic similarity representations "jump" to the last layer Xu et al. (2018b). A 256 jump connection that bypasses the GCN, directly transmitting feature proximity (textual similarity) to 257 the final embeddings, as illustrated in Figure  $\Pi$  (b). Additionally, we provide a theoretical justification 258 for this design from an information-theoretic perspective in Appendix A.1 259

260 D2: Combination of Structural and Semantic Embeddings The interaction between feature proximity and structural proximity dictates the formation of links. Previous theoretical work has demon-261 strated that local structural proximity and feature proximity are often incompatible, i.e. node pairs with 262 a large number of common neighbors tend to exhibit low feature proximity (Mao et al., 2023). Based 263 on this insight, we simplify the current baseline model in Eq 5 to  $\mathbf{Y} = MLP([\mathbf{H}^K; \beta \mathbf{H}_C + \mathbf{S}^T])$ 264 to reduce the hidden dimensions, thereby reducing time complexity during training. We proved in 265 Appendix A, that for any node pair (i, j), the approximation error of this design decreases as the 266 number of nodes increases. 267

D3: Cache Embedding Strategy Figure 1 illustrates how gradients from the learning objective of LMGJOINT are back-propagated through both the GCN and the final encoder layers of the PLM. This integration allows the PLM to capture structural features while enabling the GCN to enhance feature

16

14

12

10

70

75

Number of Params (Log Scale)

274 275 276

270 271

272

273

277 278 279

281



284 Figure 3: Illustration of the performance-complexity trade-off between LMGJOINT (green) and 285 benchmarked methods on Cora. The x-axis represents AUC, the y-axis shows the number of 286 parameters (log-scale) and the marker radius indicates inference time. Points closer to the bottom-287 right corner reflect higher cost-effectiveness.

85

PageRank

80

FT-MiniLM

90

AUC (%)

Budd

NCNC

100

105

GCN

95

Node2Vec

Proposed

GCN4LP

PLM-based

GraphEmb

○ Inference Time

Heuristic

288

298

299

300

301 302

303

305

289 aggregation. To mitigate the high memory cost of PLM training, we employ a cache embedding strategy. Node feature  $H^0$  is initialized by default node features X, we save these pre-computed 290 embeddings in a cache. In the current mini-batch, we re-encode only the tokens associated with 291 the target and source links  $(\mathbf{X}_i^T, \mathbf{X}_i^T)$  by PLM, and then concatenate them with pre-computed node 292 features as the input for GCN  $\mathbf{X} = [\mathbf{X}_{\mathcal{V} \setminus \{i,j\}}; \mathbf{X}_i^T; \mathbf{X}_j^T]$ . This approach significantly reduces the per-293 mini-batch computational cost from O(Nd) (where N is the number of nodes and d is the embedding dimension) to O(d). In summary, our method is both parameter- and memory-efficient, enabling 295 training on a single A100 GPU with 40GB VRAM across all proposed datasets. Detailed complexity 296 analysis is provided in Appendix C.1. 297

While these principles have been employed independently in previous works (Zhu et al., 2020; Wang et al. (2023), we are the first to advocate for their combined necessity. We substantiate our claims with both theoretical justifications and a comprehensive empirical analysis across diverse datasets.

#### 5 BENCHMARKING

304 In this section, we provide an overview of the various benchmarked methods. Specifically, we evaluate structure-only approaches (heuristic, structure-based) and text-only methods (PLM-Inf-MLP, 306 FT-PLM-MLP) to assess current advancements in each setting. We categorize all existing approaches 307 into two broad classes: graph-based and PLM-based. The former emphasizes improving pairwise 308 structural proximity, while the latter focuses on feature proximity.

309 310

# 5.1 GRAPH-BASED METHODS

311 This section evaluates graph-based methods which can be categorized into four groups: 1) Graph 312 **Heuristic:** Local heuristic leverages modified shared neighborhoods, including Common Neighbor 313 (CN), Adamic-Adar (AA), Resource Allocation (RA). Other global heuristics such as Katz, Shortest 314 Path and symmetric PageRank(Adamic & Adar, 2003; Page et al., 1999) consider all paths between 315 connected nodes. 2) Embedding-based Methods: It focus on proximity-preserving embedding 316 methods to model neighborhood contexts via random walks, including DeepWalk, Node2Vec and 317 LINE (Perozzi et al., 2014; Tang et al., 2015; Grover & Leskovec, 2016). 3) Aggregation-based 318 **Methods**: GCNs aggregate information recursively from first-hop neighbors (Kipf & Welling, 2016) 319 Velickovic et al., 2017). SAGE embeds self and neighboring nodes separately to handle heterophily 320 (Hamilton et al., 2017; Zhu et al., 2020). GIN achieves the same expressivity as Weisfeiler-Lehman 321 test (Xu et al., 2018a) by employing injective transformation and we perform a dot product with an MLP layer to handle the final embeddings. 4) GCN4LP: SEAL, BUDDY and ELPH introduces 322 different labeling techniques and structure embeddings to address the automorphic nodes problem 323 (Zhang et al., 2020; Zhang & Chen, 2018; Chamberlain et al., 2023). The latest GCNs augment aggregated features by leveraging local structures, such as CNs to enhance performance. Notable
examples include NCN/NCNC (Wang et al., 2023) and NeoGNN (Yun et al., 2021). Among graphbased methods, categories (1), (2) and (3) are structure-only methods since they do not utilize
node features. To sum up, we include CN, AA, RA, Katz, Shortest Path, PageRank, DeepWalk and
Node2Vec, GCN, GIN, SAGE, GAT, SEAL, NeoGNN, ELPH, BUDDY, NCNC, NCN, HLGNN.

# 330 5.2 PLM-BASED METHODS331

329

332 We introduce and transfer four extended benchmarked frameworks from node classification task, each 333 with progressively increasing resource requirements. This section provides a detailed overview of the PLMs/LLMs used and their corresponding embedding configurations. We investigate the following 334 configurations: (1) PLM/LLM as a Fixed Inference Model (PLM/LLM-Inf): The PLM/LLM 335 operates in a frozen state, generating static embeddings that are then fed into a Multi-Layer Perceptron 336 (MLP) for binary classification. Specifically, the final hidden states of the PLMs are utilized as text 337 embeddings. (2) Fine-Tuned PLMs (FT-PLM): Extending the PLM/LLM-Inf setup, fine-tuning 338 is introduced by training the last few encoder layers alongside the MLP (Chen et al.) 2024a). (3) 339 PLM/LLM-Inf-GCN: This configuration first generates text embeddings from the PLMs as node 340 features, after which a GCN is trained on the updated node features, but without further training on 341 PLM/LLM. (4) FT-PLM-GCN: Building upon (3), this approach optimizes all parameters, including 342 those in the last encoder layers of the PLM, the GCN and the MLP (Yang et al., 2021). Configurations 343 (1) and (2) are classified as text-only methods. To sum up, we benchmark cascade (3) and nested 344 architecture (4) adopted from prior work in node classification.

345 Selection of PLM/LLMs In this paper, we define PLMs as those models practical for inference 346 and fine-tuning within typical academic budgets, such as BERT (Devlin et al., 2019) and LLMs 347 refer to models requiring substantial computational resources to fine-tune such as thousands of 348 GPUs or TPUs, exemplified by Llama-3-8B (Dubey et al., 2024). We utilize both encoder-only and 349 decoder-as-embedder (BehnamGhader et al., 2024) models for text embedding, including (1) BERT, 350 a lightweight deep text embedding model pretrained in a self-supervised manner (Devlin et al., 2019); (2) e5-large-v2 (Wang et al., 2022), Sentence-Transformers MiniLM-L6-v2 (Reimers & Gurevych, 351 2019a) and MPNet(Song et al., 2020), pretrained using a contrastive learning approach. Additionally, 352 we incorporate Meta-LLaMA-3-8B (Dubey et al., 2024), a decoder-only LLM, which we include 353 as a case study for text embedding at scale. To ensure consistency and comparability with existing 354 studies, we include shallow embedding methods such as bag-of-words (Harris, 1954) and Word2Vec 355 (Mikolov et al., 2013). We leverage the [EOS] token in LLaMA3 and the [CLS] token in sentence 356 embedding models as node features and fine tune them with full-parameter tuning strategy. Further 357 details on fine-tuning and embedding strategy are provided in Appendix B.3. To sum up, we include 358 BERT, e5-large-v2, MiniLM-L6-v2, MPNet and Meta-LlaMA-3-8B, bag-of-words and Word2Vec 359 as text encoder with 4 experiment settings.

- 360 361 5.3 Experiment Settings
- 362 **Metrics choice** To ensure consistency in previous works Hu et al. (2021), we benchmark all ap-363 proaches on the Cora, PubMed and Arxiv\_2023. Furthermore, our evaluation is extended to all 364 nine datasets with metrics including Hits@50, Hits@100, MRR (Mean Reciprocal Rank) and AUC. 365 Hits@K quantifies the ratio of positive edges ranked within the top K positions, while MRR evaluates 366 the model's ability to rank positive above negative ones. AUC assesses the model's ability to score 367 positives higher than negatives, offering numerical stability and scale invariance (see Appendix D). 368 To avoid distribution shift caused by the feature-based split method (Wang et al., 2023), we apply a uniform random split of 80%, 15% and 5% across all datasets. For all experiments, the results are 369 reported on randomly sampled test edges for datasets larger than pwc\_medium in Table 3. All metrics 370 are reported as the mean and standard deviation, averaged over five random seeds. We exclude the 371 target link (link to be predicted) in each mini-batch. 372

Hyperparameter Ranges We utilize hierarchical grid search for hyperparameter optimization
 across all GCN models, tuning parameters such as learning rate, weight decay, the number of
 convolution layers, MLP layers, the number of heads in GAT and hidden neurons. For details on
 specific hyperparameters in GCN4LPs. Due to the training burden for PLM/LLMs, we use the
 same parameters for GCN and GCN4LPs in LLM-GCN related methods (i.e., PLM/LLM-Inf-GCN,
 FT-PLM-GCN and Ours) and only optimize the learning rate and weight decay. We utilize the same

378 Table 2: Benchmark results showing mean  $\pm$  stdev for Hits@50, Hits@100, and MRR metrics on 379 Cora, PubMed, and Arxiv 2023. The top 1-3 ranked models are highlighted in emerald, while ranks 380 4-6 are highlighted in green. Darker colors represent higher ranks. All results are provided under our benchmark with consistent training and evaluation settings. Heuristic, Structure, GCN, and GCN4LP 381 are existing approaches, followed by extended baselines, with the final category representing proposed 382 methods. | model | highlights the best model among the existing approaches. We integrated both soft and hard CN into our LMGJOINT. To distinguish, LMGJOINT and LMGJOINT-C refers to soft and 384 hard common neighbor introduced in Section 4 385

386											
387	Category	Models	Hits@50	Cora Hits@100	MRR	Hits@50	PubMed Hits@100	MRR	Hits@50	Arxiv 2023 Hits@100	MRR
388	Hannistia	CN	50.36±0.03	50.36±0.03	32.88±0.09	33.32±0.02	$33.32 \pm 0.02$	$21.13 \pm 0.02$	27.20±0.01	27.20±0.01	$14.66 \pm 0.06$
389	Heuristic	RA RA	50.36±0.03	50.36±0.03	47.17±0.11	$33.32\pm0.02$ $33.32\pm0.02$	33.32±0.02	23.94±0.16	27.20±0.01 27.20±0.01	27.20±0.01 27.20±0.01	19.87±0.30 19.16±0.27
390		PPR/sym	$84.74{\scriptstyle\pm0.00}$	$88.93{\scriptstyle\pm0.00}$	58.86±0.98	$69.81{\scriptstyle\pm0.02}$	$72.95{\scriptstyle\pm0.01}$	$28.04{\scriptstyle\pm0.91}$	$65.68{\scriptstyle\pm0.02}$	$67.86 {\pm 0.02}$	$26.57{\scriptstyle\pm0.82}$
204	Structure	Katz	$69.25 \pm 0.02$	$69.25 \pm 0.02$	38.17±0.12	66.02±0.02	66.02±0.02	30.94±0.08	55.39±0.01	$55.39 \pm 0.01$	21.76±0.21
391		Node2Vec	84.00±0.07 83.08±0.05	89.31±0.03 88.38±0.03	$44.39 \pm 0.96$ $39.94 \pm 1.10$	$64.33 \pm 0.01$ $63.95 \pm 0.02$	$70.95 \pm 0.01$	19.66±1.00 20.68±0.24	$30.19 \pm 0.11$	$40.81 \pm 0.20$	$4.47 \pm 0.03$ $5.60 \pm 0.03$
392		GCN	91.46±2.36	96.20±1.71	$45.84{\scriptstyle\pm8.40}$	83.11±2.19	89.59±1.73	$24.55 {\pm} 4.02$	45.07±0.87	$52.46 {\pm} 1.44$	$17.62 \pm 3.34$
393	GCNs	GAT	89.80±2.00	95.34±1.61	49.82±10.04	74.23±2.54	84.03±1.59	18.13±5.81	43.09±1.22	$53.49 \pm 1.06$	13.58±4.33
394		GIN	$91.54 \pm 2.91$	95.10±1.37 96.05±2.13	40.03±6.70 51.90±6.65	86.92±1.68	93.30±0.70 92.02±0.91	24.63±2.24	$45.42\pm3.12$ $45.35\pm2.58$	$50.09 \pm 3.13$ $53.22 \pm 2.05$	$11.32 \pm 1.67$ 14.79±4.53
395		SFAI	87 38+3.06	92 03+296	37.81+9.93	84 62+3 53	89 52+1 27	49 02+13 91	56 98+1 89	67 34+ 3 74	22 47+3 69
	CONTR	NeoGNN	81.03±3.11	90.04±2.02	41.48±5.11	73.17±5.29	81.25±8.14	31.44±3.85	64.54±11.14	$69.34 \pm 8.56$	28.07±15.62
396	GUN4LP	ELPH	$87.30 {\pm} 4.94$	$94.91 {\pm} 2.17$	$39.86{\scriptstyle\pm10.20}$	$59.19 \pm 5.58$	$74.62 \pm 1.64$	24.61±3.17	$57.66 \pm 1.55$	$66.95 {\pm}  3.62$	$29.22 \pm 5.95$
397		BUDDY	87.82±3.14	95.42±2.26	30.78±5.55	76.14±3.46	89.25±2.27	$19.46 \pm 2.42$	$52.25 \pm 2.01$	$60.49 \pm 0.94$	$18.75 \pm 3.71$
		HL-GNN	90.59±3.41	94.62±1.87	50.35±10.07	85.14±1.83	91.87±1.36	31.49±7.84	76.51±1.68	84.23±0.82	24.21±7.69
398		NCN	96.16±1.62	98.74±0.96	45.76±16.39	86.44±2.03	93.21±1.10	$25.92 \pm 4.33$	82.34±2.45	88.83± 1.43	37.92±13.21
399		NCNC	95.42±2.41	98.67±0.76	48.68±18.60	86.49±0.99	93.74±0.25	20.31±6.51	81.86±1.64	$89.13 \pm 2.08$	35.67±12.30
/00		BERT	35.79±2.50	56.90±3.26	3.42±0.47	36.12±0.37	$48.73 {\pm}~1.43$	6.56±0.70	37.66±1.57	48.74± 1.15	$10.04 \pm 0.85$
400	PLM-Inf-MLP	MiniLM	83.39±0.00	92.99±0.00	34.29±4.10	66.35±0.29	$81.90 \pm 0.03$	$21.54 \pm 0.11$ 21.70 ± 1.58	68.15±0.09	$77.62 \pm 0.03$	$16.91 \pm 0.18$
401		Llama-3-8B	$89.15\pm0.72$	$95.64\pm0.41$	$24.40\pm2.48$ 31 19+3 49	79.87±1.19	$82.39 \pm 0.26$ $89.01 \pm 0.53$	21.79±1.38 22.87+4.47	83 18+1 19	$89.91 \pm 0.24$	21.09±0.03
402		BERT	89 17+286	96 99±1 36	30.90+4.33	73 70+4 01	84 45+ 2 92	17 11+3 90	77 75+3.46	87 56+ 2.05	29 54+3 98
402		e5-large	92.09±1.70	96.92±1.35	38.63±9.39	76.26±2.55	$87.23 \pm 1.60$	19.75±5.81	80.48±2.52	89.35±1.33	31.73±6.62
403	FI-PLM-MLP	MiniLM	$92.49 \pm 2.13$	$96.68 \pm 1.69$	$35.55{\scriptstyle\pm5.82}$	$75.87 \pm 3.72$	$86.80 \pm 1.98$	$20.79 \pm 6.32$	$80.20{\scriptstyle\pm2.62}$	$88.38 {\pm} 1.06$	$29.86{\scriptstyle\pm}5.82$
404		mpnet	$93.44 \pm 1.64$	$97.78{\scriptstyle\pm0.66}$	$22.55 \pm 10.71$	63.27±31.76	90.69±2.49	9.38±3.12	$82.72 \pm 1.28$	91.44±0.75	8.42±6.49
105		MiniLM-NCNC	96.13±1.20	98.81±0.49	38.96±13.20	90.32±1.52	96.11±0.60	22.56±3.30	65.65±1.80	70.61±2.24	29.10±3.83
405	PLM-Inf-GCN	Llama-NCNC	$95.13 \pm 1.13$ $95.57 \pm 1.02$	98.81±0.74 98.73±0.65	27 45+7 86	90.80±1.95 84.65±1.95	$90.09 \pm 0.56$ 92.39 \pm 1.46	$27.02\pm 5.96$ 20 51+9 80	85.24±1.20 84.68+1.72	$90.40 \pm 1.14$ $91.90 \pm 1.33$	23.14±9.39 27.16+11.48
406		BERT-NCNC	77.47±1.77	84.11±2.70	$25.39 \pm 12.42$	72.80±1.78	82.48±1.43	$23.49 \pm 3.07$	58.83±3.91	68.50±3.49	22.80±2.55
407		MiniLM-GAT	54.23±4.08	76.99±6.58	13.74±5.21	29.44±2.84	$43.75 \pm 5.46$	4.26±1.75	13.76±2.22	25.18±4.17	2.62±0.55
100	FT-PLM-GCN	mnnet-GIN	$89.04 \pm 9.23$ $89.01 \pm 5.54$	$82.01 \pm 4.19$ 97 55+186	13.09±1.35 29.06+7.96	$45.51 \pm 3.83$ $46.82 \pm 4.22$	$63.18 \pm 1.34$ $63.42 \pm 2.77$	$11.86 \pm 3.61$	$49.11 \pm 4.22$ 55 11+3 70	$64.49 \pm 3.26$	$11.48 \pm 3.41$ 18 88+5 89
400	111201000	mpnet-GAT	74.90±9.22	86.96±6.71	23.16±11.10	35.82±3.81	$52.43 \pm 4.18$	5.36±1.69	19.50±1.91	29.43±3.60	$4.49 \pm 0.91$
409		mpnet-SAGE	$82.01{\scriptstyle\pm4.19}$	$93.88{\scriptstyle\pm0.28}$	$25.34{\scriptstyle\pm8.06}$	$57.58{\scriptstyle \pm 3.78}$	$71.62{\pm}2.78$	$11.91{\pm}3.58$	$52.97{\scriptstyle\pm}5.05$	$66.28{\scriptstyle \pm 4.50}$	$14.51{\scriptstyle \pm 3.37}$
410		MiniLM-LMGJOINT-C	99.92±0.18	99.92±0.18	41.52±19.50	99.91±0.09	99.94±0.08	44.99±10.82	90.61±2.25	98.16±1.73	35.47±10.91
/111	Ours	mpnet-LMGJOINT-C	93.28±14.16	95.81±9.15	28.92±7.14	99.27±1.19	$99.22 \pm 4.96$ $99.95 \pm 0.08$	$22.72 \pm 1.51$ 23.99 $\pm 11.63$	73.09±16.32	77.99±17.20	$14.68 \pm 6.17$
411		mpnet-LMGJOINT	$100.00{\scriptstyle\pm0.00}$	$100.00{\scriptstyle\pm0.00}$	$68.43 \pm 14.23$	91.67±4.96	$97.13 \pm 1.74$	31.66±5.33	$89.17 \pm 5.45$	$94.85 \pm 3.15$	45.70±3.88
412		e5-large-LMGJOINT-C	$99.92{\scriptstyle\pm0.18}$	$99.92{\scriptstyle\pm0.18}$	$41.46 {\pm} 25.49$	99.11±1.54	$100.00 \pm 0.00$	21.66±9.66	83.97±4.23	$98.01 \pm 0.72$	$12.66 \pm 2.66$
		e5-large-LMGJOINT	$96.29 \pm 2.08$	$98.89 \pm 1.02$	65.26±11.52	77.34±2.19	$88.41 \pm 1.14$	23.80±3.29	$80.01 \pm 2.53$	87.71±1.48	$42.02 \pm 5.56$

413

414

416 417 418

419 420

415

parameters of GCN and LLM in our proposed method. Further details about parameter tuning and experiment setting are provided in Appendix B

#### **BENCHMARK ANALYSIS** 6

421 We analyze the benchmark results by addressing the following questions: (1) Is utilizing PLM alone 422 more effective than a structure-based method? (2) Does the previous GCN4LP SOTA persist under the new configurations? (3) Should PLMs and GCN-based methods be trained separately? 423

424 **Text-only vs. Structure-only.** To address (1), we compare the performance between structure-425 only methods (Heuristic, Structure) and text-only approaches (PLM/LLM-Inf-MLP, FT-PLM-MLP) 426 in Table 2. Although structure-only methods achieve better performance on Cora and PubMed 427 when assessing MRR, text-only methods substantially show better performance in other metrics 428 including Hits@50 and Hits@100. Furthermore, while assessing AUC as shown in Figure 3, the performance of FT-MiniLM, a text-only model, approaches the similar performance of a robust 429 GCN4LP method NCNC. This suggests that PLM-based models can achieve strong performance even 430 in the absence of topological information assessed Hits@K and AUC. However, it fails to outperform 431 the structure-based method in MRR across benchmarked datasets under a unified experiment setting.

Table 3: Results on extensive datasets: Comparison with the strongest baseline in each category using
 AUC. Mean accuracy ± standard deviation is reported across different data splits. The best model for
 each benchmark is highlighted in emerald.

	Small			Medium				LARGE	
	Pwc <sub>small</sub>	Cora	Arxiv <sub>2023</sub>	PubMed	Pwc <sub>medium</sub>	History	Photo	Ogbn-arxiv	Citation
Embedding-MLP: No	on-contextualize	ed Shallow Em	oeddings						
TF-IDF	$63.50 \pm 8.59$	$68.27 \pm 2.52$	$76.65 \pm 1.95$	$67.06 \pm 2.34$	$70.94 \pm 0.94$	$60.94 \pm 0.48$	$62.80 \pm 0.62$	$62.63 \pm 2.88$	$57.18 \pm$
Word2Vec	$51.00 \pm 2.24$	$60.15 \pm 1.77$	$85.22\pm0.92$	$83.88 \pm 0.91$	$81.79 \pm 0.55$	$65.54 \pm 0.21$	$64.71 \pm 0.12$	$85.57\pm0.28$	$80.50 \pm$
PLM-Inf-MLP: Loca	I Sentence Emb	edding Models							
MiniLM-L6-v2	$51.90 \pm 11.54$	$91.2\bar{2} \pm 0.04$	$95.22\pm0.00$	$96.20\pm0.00$	$98.39 \pm 0.04$	$94.64 \pm 0.01$	$84.14\pm0.03$	$98.22\pm0.01$	$97.86 \pm 0$
e5-large-v2	$80.60 \pm 2.57$	$83.87 \pm 0.23$	$95.72 \pm 0.01$	$96.73 \pm 0.03$	$97.83 \pm 0.01$	$95.97 \pm 0.00$	$85.04 \pm 0.44$	$97.92 \pm 0.01$	$98.05 \pm 0$
BERT	$69.85 \pm 2.40$	$65.09 \pm 1.41$	$86.37 \pm 0.27$	$88.96 \pm 0.31$	$83.85\pm0.33$	$90.26 \pm 0.36$	$73.12\pm0.76$	$86.89 \pm 0.26$	$86.22 \pm$
LLM-Inf-MLP									
Llama-3-8B	$94.65 \pm 1.23$	$92.60\pm0.12$	$97.62\pm0.03$	$98.09 \pm 0.10$	$97.74 \pm 0.03$	$97.28 \pm 0.08$	$88.62\pm0.26$	$99.06 \pm 0.05$	$99.05 \pm$
Strong GCN4LP base	eline								
NCN	$86.65 \pm 5.37$	$96.66 \pm 1.14$	$97.30 \pm 0.26$	$98.66 \pm 0.18$	$98.46 \pm 0.19$	$97.77 \pm 0.30$	$96.58 \pm 0.2$	$98.96 \pm 0.07$	$98.18 \pm 0$
NCNC	$86.87 \pm 7.99$	$96.56 \pm 1.04$	$97.42 \pm 0.37$	$98.66 \pm 0.12$	$98.45 \pm 0.21$	$97.79 \pm 0.25$	$96.79 \pm 0.25$	$98.93 \pm 0.13$	$98.68 \pm$
FT-PLM-MLP									
mpnet-FT	$85.93 \pm 5.86$	$94.71 \pm 1.16$	$97.36 \pm 0.33$	$98.06 \pm 0.19$	$97.72 \pm 0.54$	$93.91 \pm 0.64$	$82.26 \pm 0.92$	$98.21 \pm 0.26$	$98.17 \pm$
e5-large-v2-FT	$86.95 \pm 4.93$	$94.27 \pm 0.85$	$97.39 \pm 0.33$	$97.64 \pm 0.36$	$94.06 \pm 0.84$	$94.82 \pm 1.18$	$86.26 \pm 1.16$	$97.68 \pm 0.17$	$97.08 \pm 3$
MiniLM-L6-v2-FT	$87.09 \pm 2.51$	$93.98 \pm 0.85$	$97.25 \pm 0.36$	$97.79 \pm 0.14$	$97.95 \pm 0.44$	$95.58 \pm 0.51$	$88.34 \pm 0.75$	$98.80 \pm 0.16$	$97.85 \pm$
PLM-Inf-GCN									
MiniLM-NCN	$87.15 \pm 6.84$	$96.93 \pm 0.54$	$86.69 \pm 0.28$	$98.97 \pm 0.10$	$98.99 \pm 0.16$	$99.42 \pm 0.11$	$99.59 \pm 0.03$	$99.58 \pm 0.07$	$98.17 \pm$
e5-large-NCN	$88.31 \pm 4.99$	$96.72\pm0.67$	$97.82 \pm 0.22$	$97.24 \pm 0.20$	$99.03 \pm 0.18$	$99.50 \pm 0.13$	$99.56 \pm 0.04$	$99.52\pm0.07$	$98.15 \pm$
Ours									
MiniLM-LMGJOINT	$88.20 \pm 5.93$	$97.79 \pm 0.66$	$98.22 \pm 0.30$	$98.30 \pm 0.51$	$99.00 \pm 0.15$	99.14 + 0.02	$99.58 \pm 0.04$	$99.60 \pm 0.07$	$99.54 \pm$
mpnet-LMGJOINT	$89.36 \pm 5.37$	$98.78 \pm 1.02$	$98.79 \pm 0.49$	$99.34 \pm 0.22$	$99.34 \pm 0.09$	$99.54 \pm 0.01$	$99.63 \pm 0.02$	$99.72 \pm 0.04$	$99.76 \pm$

Does GCN4LP Maintain Its Superiority? Currently promising GCN4LP methods—such as SEAL, BUDDY and NeoGNN, do **NOT** show a significant advantage over GCNs in this setting. Nonetheless, NCN/NCNC maintains superior performance across both Hits@k and MRR metrics, establishing itself as the strongest baseline with optimal computational complexity, as shown in Figure 3 Overall, all aggregated methods, including both GCN and GCN4LP, consistently outperform the structure-only methods, reaffirming the effectiveness of GCN-based approaches as a robust foundational framework. In PLM-Inf-GCN category, we observed that NCNC's performance could be improved simply by replacing the original node feature with PLM-based text embeddings. This signifies the significant potential of PLM-based text embeddings to enhance performance in link prediction tasks.

Cascade PLM-GCN vs. Nested PLM-GCN. To answer Q3, we evaluate the impact of fine-tuning within cascade and nested frameworks by comparing NCNC and PLM-Inf-NCNC, PLM-Inf-MLP and FT-PLM-MLP, GCN and FT-PLM-GCN. We observe limited improvement in Hits@K, accompanied by a notable decline in MRR. It indicates that incorporating context-aware PLM text encoding can enhance representation quality concerning Hits@K. However, the performance decay in MRR is consistent with prior findings by Chen et al. (2024b), which suggest that fine-tuning without specific design considerations can sometimes result in negative performance gains. In summary, these results indicate that optimizing PLMs independently of topology-based methods does not lead to consistent improvements across all metrics. This underscores the importance of developing an effective nested architecture to fully leverage the strengths of both approaches. 

# 6.1 EMPIRICAL EVALUATION OF THE PROPOSED METHOD

We evaluate LMGJOINT through two complementary studies. First, a horizontal analysis in Experiment 1 (Exp) compares LMGJOINT against diverse graph-based methods on three popular datasets
(Table 2). Second, a vertical analysis in Exp 2 examines competitive baselines in Exp 1 and broader
LM-based approaches, extending the evaluation to nine datasets (Table 3) to assess the generality and robustness of LMGJOINT's improvements.

Exp 1: Horizontal perspective From Table 2. LMGJOINT consistently demonstrates superior performance across the three datasets and extensive metrics, outperforming the second-best category, PLM-Inf-GCN and strong baselines such as NCNC and LLaMA3-inf-MLP. It achieves the best results in 7 out of 9 comparative evaluations, highlighting its robustness and effectiveness. Compared to all benchmarked models, PLM/LLM-Inf-GCN family excels in Hits@K predominantly, while those with the highest MRR scores are concentrated within the Structure and GCN4LP categories. In contrast, our approach excels in both Hits@K and MRR, indicating its ability to preserve both the feature proximity from sentence embeddings and the structural proximity from graph topology.

486 Exp 2: Vertical perspective In Table 3, we selected powerful baselines from Exp 1, includ-487 ing NCN(C), FT-PLM and PLM/LLM-Inf-NCN and extend to all proposed graphs (Except ogbn-488 paper100M due to GPU limitation). When comparing LM/LLM from a vertical perspective, 489 Word2Vec outperforms the bag-of-words method among non-contextual embeddings, while lo-490 cal sentence embeddings improve AUC performance by over 10% compared to non-contextual embeddings in AUC. We observe a strong positive correlation between model performance gains and 491 the number of parameters in PLM. LLaMA3 leads among LMs, achieving the best performance on 492 pwc\_small, indicating that text-only methods with large PLMs excel when structural information is 493 limited. Nevertheless, Our LMGJOINT consistently outperforms these baselines across all evaluated 494 datasets, achieving up to a 2.6% improvement in AUC. 495

496 497

498

7 CONCLUSION

499 This work tackles the under-explored area of joint PLM-GCN architecture design for link prediction 500 by benchmarking PLM and GCN-based methods on extensive datasets. In our benchmark, we focused on discussing the impact of fine-tuning modules and text embeddings within various PLM-GCN 501 architectures. We introduce the LMGJOINT, a simple yet powerful nested framework that combines 502 the strengths of both GCN4LP and PLM-based methods while avoiding their weaknesses. Extensive 503 and rigorous experiments demonstrate that LMGJOINT consistently improves performance across all 504 the metrics on various datasets, with minimal hyperparameter tuning required. We expect it to benefit 505 PLM-GCN models and other graph learning tasks, including node classification and regression. 506

507 508

509

510

511

512

529

530

531

# References

- Lada A Adamic and Eytan Adar. Friends and neighbors on the web. *Social Networks*, 25(3):211–230, 2003. ISSN 0378-8733. doi: https://doi.org/10.1016/S0378-8733(03)00009-1. URL https://www.sciencedirect.com/science/article/pii/S0378873303000091.
- Parishad BehnamGhader, Vaibhav Adlakha, Marius Mosbach, Dzmitry Bahdanau, Nicolas Chapados, and Siva Reddy. Llm2vec: Large language models are secretly powerful text encoders. *ArXiv*, abs/2404.05961, 2024. URL https://api.semanticscholar.org/CorpusID:
  269009682.
- Markus Brede. Networks—an introduction. mark e. j. newman. (2010, oxford university press.)
  \$65.38, £35.96 (hardcover), 772 pages. isbn-978-0-19-920665-0. Artificial Life, 18:241-242, 2012.
  URL https://api.semanticscholar.org/CorpusID:207677121
- Benjamin Paul Chamberlain, Sergey Shirobokov, Emanuele Rossi, Fabrizio Frasca, Thomas Markovich, Nils Hammerla, Michael M. Bronstein, and Max Hansmire. Graph Neural Networks for Link Prediction with Subgraph Sketching, May 2023. URL http://arxiv.org/ abs/2209.15486 arXiv:2209.15486 [cs].
- Ines Chami, Rex Ying, Christopher Ré, and Jure Leskovec. Hyperbolic graph convolutional neural networks. Advances in neural information processing systems, 32:4869–4880, 2019. URL https://api.semanticscholar.org/CorpusID:202784587.
  - Zhikai Chen, Haitao Mao, Hang Li, Wei Jin, Hongzhi Wen, Xiaochi Wei, Shuaiqiang Wang, Dawei Yin, Wenqi Fan, Hui Liu, et al. Exploring the potential of large language models (llms) in learning on graphs. *ACM SIGKDD Explorations Newsletter*, 25(2):42–61, 2024a.
- Zhikai Chen, Haitao Mao, Jingzhe Liu, Yu Song, Bingheng Li, Wei dong Jin, Bahare Fatemi, Anton Tsitsulin, Bryan Perozzi, Hui Liu, and Jiliang Tang. Text-space graph foundation models: Comprehensive benchmarks and new insights. *ArXiv*, abs/2406.10727, 2024b. URL https: //api.semanticscholar.org/CorpusID:270559362.
- 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. *ArXiv*, abs/2111.00064, 2021. URL https://api.semanticscholar.org/CorpusID: 240354406.

540 541 542 543	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Development)</i>
544 545	Short Papers), pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL https://aclanthology.org/N19-1423.
546 547 548 549	Keyu Duan, Qian Liu, Tat-Seng Chua, Shuicheng YAN, Wei Tsang Ooi, Michael Qizhe Xie, and Junxian He. Simteg: A frustratingly simple approach improves textual graph learning, 2024. URL <a href="https://openreview.net/forum?id=EFGwiZ2pAW">https://openreview.net/forum?id=EFGwiZ2pAW</a> .
550 551 552	Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. <i>arXiv preprint arXiv:2407.21783</i> , 2024.
555 555 556 557 558	Aditya Grover and Jure Leskovec. node2vec: Scalable Feature Learning for Networks. In <i>Proceedings</i> of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 855–864, San Francisco California USA, August 2016. ACM. ISBN 978-1-4503-4232-2. doi: 10.1145/2939672.2939754. URL https://dl.acm.org/doi/10.1145/2939672.2939754. 2939754.
559 560 561 562	William L. Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. In <i>Neural Information Processing Systems</i> , 2017. URL https://api.semanticscholar.org/CorpusID:4755450.
563	ZS Harris. Distributional structure. Word, pp. 146-162, 1954.
564 565 566 567	Xiaoxin He, Xavier Bresson, Thomas Laurent, Adam Perold, Yann LeCun, and Bryan Hooi. Harness- ing explanations: Llm-to-lm interpreter for enhanced text-attributed graph representation learning, 2023.
568 569 570	Zhongmou He, Jing Zhu, Shengyi Qian, Joyce Chai, and Danai Koutra. Linkgpt: Teaching large language models to predict missing links. <i>ArXiv</i> , abs/2406.04640, 2024. URL https://api.semanticscholar.org/CorpusID:270357491.
572 573 574	Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs. In <i>Advances in neural information processing systems</i> , volume 33, pp. 22118–22133, 2020a.
575 576 577 578	Weihua Hu, Matthias Fey, Hongyu Ren, Maho Nakata, Yuxiao Dong, and Jure Leskovec. Ogb-lsc: A large-scale challenge for machine learning on graphs. <i>ArXiv</i> , abs/2103.09430, 2021. URL https://api.semanticscholar.org/CorpusID:232257683.
579 580 581 582	Ziniu Hu, Yuxiao Dong, Kuansan Wang, and Yizhou Sun. Heterogeneous graph transformer. In <i>Proceedings of The Web Conference 2020</i> , WWW '20, pp. 2704–2710, New York, NY, USA, 2020b. Association for Computing Machinery. ISBN 9781450370233. doi: 10.1145/3366423.3380027. URL https://doi.org/10.1145/3366423.3380027.
583 584 585 586 587	Zan Huang, Xin Li, and Hsinchun Chen. Link prediction approach to collaborative filtering. In <i>Proceedings of the 5th ACM/IEEE-CS Joint Conference on Digital Libraries</i> , JCDL '05, pp. 141–142, New York, NY, USA, 2005. Association for Computing Machinery. ISBN 1581138768. doi: 10.1145/1065385.1065415. URL https://doi.org/10.1145/1065385.1065415.
588 589 590 591	Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. <i>CoRR</i> , abs/1412.6980, 2014. URL https://api.semanticscholar.org/CorpusID: 6628106.
592 593	Thomas Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. <i>ArXiv</i> , abs/1609.02907, 2016. URL https://api.semanticscholar.org/CorpusID: 3144218.

594 595 596 597	Meng-Chieh Lee, Haiyang Yu, Jian Zhang, Vassilis N. Ioannidis, Xiang song, Soji Adeshina, Da Zheng, and Christos Faloutsos. Netinfof framework: Measuring and exploiting network usable information. In <i>The Twelfth International Conference on Learning Representations</i> , 2024. URL https://openreview.net/forum?id=KY8ZNcljVU.
598 599 600 601 602	Juanhui Li, Harry Shomer, Haitao Mao, Shenglai Zeng, Yao Ma, Neil Shah, Jiliang Tang, and Dawei Yin. Evaluating graph neural networks for link prediction: Current pitfalls and new benchmarking. <i>ArXiv</i> , abs/2306.10453, 2023. URL https://api.semanticscholar.org/CorpusID: 259204112.
603 604 605 606	Zhuofeng Li, Zixing Gou, Xiangnan Zhang, Zhongyuan Liu, Sirui Li, Yuntong Hu, Chen Ling, Zheng Zhang, and Liang Zhao. Teg-db: A comprehensive dataset and benchmark of textual-edge graphs, 2024.
607	Qi Liu, Maximilian Nickel, and Douwe Kiela. Hyperbolic Graph Neural Networks. pp. 12.
608 609 610 611 612	Sitao Luan, Chenqing Hua, Minkai Xu, Qincheng Lu, Jiaqi Zhu, Xiao-Wen Chang, Jie Fu, Jure Leskovec, and Doina Precup. When Do Graph Neural Networks Help with Node Classification? Investigating the Impact of Homophily Principle on Node Distinguishability. <i>arXiv e-prints</i> , art. arXiv:2304.14274, April 2023. doi: 10.48550/arXiv.2304.14274.
613 614 615 616 617	Sitao Luan, Chenqing Hua, Minkai Xu, Qincheng Lu, Jiaqi Zhu, Xiao-Wen Chang, Jie Fu, Jure Leskovec, and Doina Precup. When do graph neural networks help with node classification? investigating the homophily principle on node distinguishability. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> , 2023. URL <a href="https://openreview.net/forum?lid=kJmYu3Ti2z">https://openreview.net/forum?lid=kJmYu3Ti2z</a>
618 619 620 621	Haitao Mao, Juanhui Li, Harry Shomer, Bingheng Li, Wenqi Fan, Yao Ma, Tong Zhao, Neil Shah, and Jiliang Tang. Revisiting link prediction: A data perspective. <i>arXiv preprint arXiv:2310.00793</i> , 2023.
622 623 624 625	Andrew Kachites McCallum, Kamal Nigam, Jason Rennie, and Kristie Seymore. Automating the construction of internet portals with machine learning. <i>Information Retrieval</i> , 3(2):127–163, 2000. ISSN 1573-7659. doi: 10.1023/A:1009953814988. URL https://doi.org/10.1023/A: 1009953814988.
626 627 628 629	Tomas Mikolov, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. In <i>International Conference on Learning Representations</i> , 2013. URL https://api.semanticscholar.org/CorpusID:5959482.
630 631 632	Mark E. J. Newman. Clustering and preferential attachment in growing networks. <i>Physical review</i> . <i>E, Statistical, nonlinear, and soft matter physics</i> , 64 2 Pt 2:025102, 2001. URL <a href="https://api.semanticscholar.org/CorpusID:9744376">https://api.semanticscholar.org/CorpusID:9744376</a> .
633 634 635 636	Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The pagerank citation ranking: Bringing order to the web. Technical Report 1999-66, Stanford InfoLab, November 1999. URL http://ilpubs.stanford.edu:8090/422/, Previous number = SIDL-WP-1999-0120.
637 638 639 640	Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. DeepWalk: Online Learning of Social Representations. In <i>Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining</i> , pp. 701–710, August 2014. doi: 10.1145/2623330.2623732. URL http://arxiv.org/abs/1403.6652, arXiv:1403.6652 [cs].
641 642 643 644	Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In Conference on Empirical Methods in Natural Language Processing, 2019a. URL https: //api.semanticscholar.org/CorpusID:201646309.
645 646 647	Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing</i> . Association for Computational Linguistics, 11 2019b. URL <a href="http://arxiv.org/abs/1908">http://arxiv.org/abs/1908</a> , 10084

648 649 650	Tarek Saier, Youxiang Dong, and Michael Färber. Cocon: A data set on combined contextualized research artifact use. 2023 ACM/IEEE Joint Conference on Digital Libraries (JCDL), pp. 47–50, 2023. URL https://api.semanticscholar.org/CorpusID:257766547.
652 653	Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. Collective classification in network data. <i>AI magazine</i> , 29(3):93–93, 2008a.
654 655 656 657 658	Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. Collective classification in network data. <i>AI Magazine</i> , 29(3):93, Sep. 2008b. doi: 10.1609/aimag. v29i3.2157. URL https://ojs.aaai.org/aimagazine/index.php/aimagazine/ article/view/2157.
659 660	Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann. Pitfalls of graph neural network evaluation. <i>arXiv preprint arXiv:1811.05868</i> , 2018.
661 662 663	Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. Mpnet: Masked and permuted pre-training for language understanding. <i>ArXiv</i> , abs/2004.09297, 2020. URL https://api.semanticscholar.org/CorpusID:215827489.
665 666 667 668 669 670	Damian Szklarczyk, Annika L Gable, David Lyon, Alexander Junge, Stefan Wyder, Jaime Huerta- Cepas, Milan Simonovic, Nadezhda T Doncheva, John H Morris, Peer Bork, Lars J Jensen, and Christianvon Mering. STRING v11: protein–protein association networks with increased coverage, supporting functional discovery in genome-wide experimental datasets. <i>Nucleic Acids</i> <i>Research</i> , 47(D1):D607–D613, 11 2018. ISSN 0305-1048. doi: 10.1093/nar/gky1131. URL https://doi.org/10.1093/nar/gky1131.
671 672 673	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. <i>ArXiv</i> , abs/2310.13023, 2023. URL https://api.semanticscholar.org/CorpusID:264405943.
674 675 676 677 678	Jian Tang, Meng Qu, Mingzhe Wang, Ming Zhang, Jun Yan, and Qiaozhu Mei. LINE: Large- scale Information Network Embedding. In <i>Proceedings of the 24th International Conference</i> <i>on World Wide Web</i> , pp. 1067–1077, May 2015. doi: 10.1145/2736277.2741093. URL http: //arxiv.org/abs/1503.03578. arXiv:1503.03578 [cs].
679 680 681 682 683	Jie Tang, Jing Zhang, Limin Yao, Juanzi Li, Li Zhang, and Zhong Su. Arnetminer: extraction and mining of academic social networks. In <i>Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining</i> , KDD '08, pp. 990–998, New York, NY, USA, 2008. Association for Computing Machinery. ISBN 9781605581934. doi: 10.1145/1401890. 1402008. URL https://doi.org/10.1145/1401890.1402008.
684 685 686	Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio', and Yoshua Bengio. Graph attention networks. <i>ArXiv</i> , abs/1710.10903, 2017. URL https://api. semanticscholar.org/CorpusID:3292002.
687 688 689 690	Kuansan Wang, Zhihong Shen, Chiyuan Huang, Chieh-Han Wu, Yuxiao Dong, and Anshul Kanakia. Microsoft academic graph: When experts are not enough. <i>Quantitative Science Studies</i> , 1(1): 396–413, 2020.
691 692 693 694	Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. Text embeddings by weakly-supervised contrastive pre-training. <i>ArXiv</i> , abs/2212.03533, 2022. URL https://api.semanticscholar.org/CorpusID: 254366618
695 696 697 698	Xiyuan Wang, Haotong Yang, and Muhan Zhang. Neural Common Neighbor with Com- pletion for Link Prediction, April 2023. URL <a href="http://arxiv.org/abs/2302.00890">http://arxiv.org/abs/2302.00890</a> , arXiv:2302.00890 [cs].
699 700 701	Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and Philip S. Yu. A Comprehensive Survey on Graph Neural Networks. <i>IEEE Transactions on Neural Networks and Learning Systems</i> , 32(1):4–24, January 2021. ISSN 2162-237X, 2162-2388. doi: 10.1109/TNNLS. 2020.2978386. URL http://arxiv.org/abs/1901.00596] arXiv:1901.00596 [cs, stat].