# SimTeG: A Frustratingly <u>Simple Approach Improves Textual Graph</u> Learning

**Anonymous ACL submission** 

#### Abstract

Textual graphs (TGs) are graphs whose nodes 001 correspond to text (sentences or documents), which are widely prevalent. The representation learning of TGs involves two stages: (i)unsupervised feature extraction and (ii) super-006 vised graph representation learning. In recent years, extensive efforts have been devoted to the latter stage, where Graph Neural Networks (GNNs) have dominated. However, the former stage for most existing graph benchmarks still relies on traditional feature engineering techniques. This motivates us to investigate the outcomes of enhancing only the text embeddings in benchmark models. While it is anticipated that advanced text embeddings will 016 boost GNN performance, key questions remain underexplored: the extent of this improvement, 017 018 particularly how advanced text features can enhance a rudimentary GNN architecture. Therefore, in this work, we investigate the impact of enhancing benchmark text embeddings exclusively and evaluate it on two fundamental graph 022 representation learning tasks: node classification and link prediction. Through extensive 024 experiments, we show that better text embeddings significantly improves the performance of various GNNs, especially basic GNN base-028 lines, on multiple graph benchmarks. Remarkably, when additional supporting text provided by large language models (LLMs) is included, a simple two-layer GraphSAGE trained on an ensemble of text embeddings achieves an accuracy of 77.48% on OGBN-Arxiv, comparable to state-of-the-art (SOTA) performance obtained from far more complicated GNN architectures. We will release our code and generated node 037 features soon.

#### 1 Introduction

040

041

042

Textual Graphs (TGs) offer a graph-based representation of text data where relationships between phrases, sentences, or documents are depicted through edges. TGs are ubiquitous in real-world applications, including citation graphs (Hu et al., 2020; Yang et al., 2016), knowledge graphs (Wang et al., 2021), and social networks (Zeng et al., 2019; Hamilton et al., 2017), provided that each entity can be represented as text. Different from traditional NLP tasks, instances in TGs are correlated with each other, which provides non-trivial and specific information for downstream tasks. In general, graph benchmarks are usually task-specific (Hu et al., 2020), and most TGs are designed for two fundamental tasks: *node classification* and *link prediction*. For the first one, we aim to predict the category of unlabeled nodes while for the second one, our goal is to predict missing links among nodes. For both tasks, text attributes offer critical information.

045

047

048

051

052

053

054

058

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

081

In recent years, TG representation learning follows a two-stage paradigm: (i) upstream: unsupervised feature extraction that encodes text into numeric embeddings, and (ii) downstream: supervised graph representation learning that further transform the embeddings utilizing the graph structure. While Graph Neural Networks (GNNs) have dominated the latter stage, with an extensive body of academic research published, the former stage surprisingly still relies on traditional feature engineering techniques. For example, in most existing graph benchmarks (Hu et al., 2020; Yang et al., 2016; Zeng et al., 2019), node features are constructed using skip-gram (Mikolov et al., 2013). This intuitively limits the performance of downstream GNNs, as it fails to fully capture textual semantics, fostering an increasing number of GNN models with more and more complex structures.

Consequently, our research aims to investigate the impact of enhancing benchmark text embeddings exclusively. While an anticipated outcome is an improved GNN performance by introducing advanced text embeddings, key inquiries remain underexplored: the extent of this improvement, and specifically, the potential enhancement of a basic GNN architecture by advanced text features. This

inquiry holds practical significance, as industry applications of GNN architectures are limited by com-086 putational efficiency. To date, a notable exception 087 is Pinsage (Ying et al., 2018), a GraphSAGE-based recommendation system for Pinterest. If incorporating advanced text features could bypass the ne-090 cessity of using complex GNN models, it would significantly boost the application of GNNs in industry. We take an step forwards to explore the above research questions by introducing a simple 094 and straightforward framework SimTeG on TGs and empirically evaluating it on two fundamental graph tasks: node classification and link prediction. We first parameter-efficiently finetune (PEFT) an LM on the textual corpus of a TG with task-specific labels and then use the finetuned LM to generate node representations given its text by removing the 101 head layer. Afterward, a GNN is trained with the 102 derived node embeddings on the same downstream 103 task for final evaluation. with extensive experi-104 ments on three prestigious graph benchmarks on 105 node classification and link prediction, we find several key observations:

• Good language modeling could generally improve the learning of GNNs on both node classification and link prediction. We evaluate SimTeG on three prestigious graph benchmarks for either node classification or link prediction, and find that SimTeG consistently outperforms the official features and the features generated by pretrained LMs (without finetuning) by a large margin. Notably, backed with SOTA GNN, we achieve *new SOTA performance* of 78.02% on OGBN-Arxiv. See Sec. 5.1 and Appendix A1 for details.

109

110

111

112

113

114

115

116

117

118

119

120

121

122

124

125

126

127

129

130

131

132

133

134

135

Incorporating advanced text features, a simple two-large GraphSAGE achieves on-par SOTA performance on node classification and link prediction tasks. Notably, a simple two-layer GraphSAGE (Hamilton et al., 2017) trained on SimTeG with proper LM backbones achieves on-par SOTA performance of 77.48% on OGBN-Arxiv (Hu et al., 2020). To date, It achieves the top three rank on the leaderboard, while the original result for sole GraphSAGE is ranked 62.

• PEFT are crucial when finetuning LMs to generate representative embeddings, because fullfinetuning usually leads to extreme overfitting due to its large parameter space and the caused fitting ability. The overfitting in the LM finetuning stage will hinder the training of downstream GNNs with a collapsed feature space. See Sec. 5.2 for details. • SimTeG is moderately sensitive to the selection of LMs. Generally, the performance of SimTeG is positively correlated with the corresponding LM's performance on text embedding tasks, e.g. classification and retrieval. In addition, the performance is not closely correlated with the number of parameters in the LM. We refer to Sec. 5.3 for details. Based on this, we expect further improvement of SimTeG once more powerful LMs for text embedding are available.

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

#### 2 Related Works

Leveraging LMs on TGs. Focusing on leveraging the power of LMs to TGs, there are several works that are existed and directly comparable with ours. For these works, they either focus on (i) designing specific strategies to generate node embeddings using LMs (He et al., 2023; Chien et al., 2021) or (ii) jointly training LMs and GNNs within a framework (Zhao et al., 2022; Mavromatis et al., 2023). Representatively, for the former one, Chien et al. (2021) proposed a self-supervised graph learning task integrating XR-Transformers (Zhang et al., 2021b) to extract node representation, which shows superior performance on multiple graph benchmarks, validating the necessity for acquiring highquality node features for attributed graphs. Jin et al. (2023) proposed two pretraining strategies for network-contextualized masked language modeling and masked node prediction to capture semantics and structure information at once. Besides, He et al. (2023) utilizes ChatGPT (OpenAI, 2023) to generate additional supporting text with LLMs. For the latter mechanism, Zhao et al. (2022) proposed a variational expectation maximization joint-training framework for LMs and GNNs to learn powerful graph representations. Mavromatis et al. (2023) designs a graph structure-aware framework to distill the knowledge from GNNs to LMs. Generally, the joint-training framework requires specific communication between LMs and GNNs, e.g. pseudo labels (Zhao et al., 2022) or hidden states (Mavromatis et al., 2023). It is worth noting that the concurrent work He et al. (2023) proposed a close method to ours. However, He et al. (2023) focuses on generating additional informative texts for nodes with LLMs, which is specifically for citation networks on node classification task. In contrast, we focus on generally investigating the effectiveness of our proposed method, which could be widely applied to unlimited datasets and tasks. Utilizing the additional text provided by He et al.

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

259

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

(2023), we further show that our method could achieve now SOTA on OGBN-Arxiv. In addition to the main streams, there are related works trying to fuse the architecture of LM and GNN for end-toend training. Yang et al. (2021) proposed a nested architecture by injecting GNN layers into LM layers. However, due to the natural incompatibleness regarding training batch size, this architecture only allows 1-hop message passing, which significantly reduce the learning capability of GNNs.

187

188

189

190

191

192

193

194

195

196

197

201

207

208

209

210

211

212

213

214

215

216

217

219

221

229

234

More "Related" Works. O Graph Transformers (Wu et al., 2021; Ying et al., 2021; Hussain et al., 2022; Park et al., 2022; Chen et al., 2022): Nowadays, Graph Transformers are mostly used to denote Transformer-based architectures that embed both topological structure and node features. Different from our work, these models focus on graphlevel problems (e.g. graph classification and graph generation) and specific domains (e.g. molecular datasets and protein association networks), which cannot be adopted on TGs. 2 Leveraging GNNs on Texts (Zhu et al., 2021; Huang et al., 2019; Zhang et al., 2020): Another seemingly related line on integrating GNNs and LMs is conversely applying GNNs to textual documents. Different from TGs, GNNs here do not rely on ground-truth graph structures but the self-constructed or synthetic ones.

#### 3 Preliminaries

Notations. To make notations consistent, we use **bold** uppercase letters to denote matrices and vectors, and calligraphic font types (e.g.  $\mathcal{T}$ ) to denote sets. We denote a textual graph as a set  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{T})$ , where  $\mathcal{V}$  and  $\mathcal{E}$  are a set of nodes and edges, respectively.  ${\mathcal T}$  is a set of text and each textual item is aligned with a node  $v \in \mathcal{V}$ . For practical usage, we usually rewrite  $\mathcal{E}$  into  $\mathbf{A} \in$  $\{0,1\}^{|\mathcal{V}|\times|\mathcal{V}|}$ , which is a sparse matrix, where entry  $\mathbf{A}_{i,j}$  denotes the link between node  $v_i, v_j \in \mathcal{V}$ . Problem Formulations. We focus on two fundamental tasks in TGs: (i) node classification and (ii) link prediction. For node classification, given a TG  $\mathcal{G}$ , we aim to learn a model  $\Phi: \mathcal{V} \to \mathcal{Y}$ , where  $\mathcal{Y}$  is the ground truth labels. For link prediction, given a TG  $\mathcal{G}$ , we aim to learn a model  $\Phi : \mathcal{V} \times \mathcal{V} \to \{0, 1\},\$ where  $f(v_i, v_j) = 1$  if there is a link between  $v_i$ and  $v_i$ , otherwise  $f(v_i, v_j) = 0$ . Different from traditional tasks that are widely explored by the graph learning community, evolving original text into learning is non-trivial. Particularly, when ablating the graphs structure, node classification and

link prediction problem are collapsed to text classification and text similarity problem, respectively. This sheds light on how to leverage LMs for TG representation learning.

Node-level Graph Neural Networks. Nowadays, GNNs have dominated graph-related tasks. Here we focus on GNN models working on node-level tasks (i.e. node classification and link prediction). These models work on generating node representations by recursively aggregating features from their multi-hop neighbors, which is usually noted as message passing. Generally, one can formulate a graph convolution layer as:  $X_{l+1} = \Psi_l(CX_l)$ , where C is the graph convolution matrix (e.g.  $C = D^{-1/2}AD^{-1/2}$  in Vanilla GCN (Kipf and Welling, 2016)) and  $\Psi_l$  is the feature transformation matrix. For the node classification problem, a classifier (e.g., an MLP) is usually appended to the output of a k-layer GNN model; while for link prediction, a similarity function is applied to the final output to compute the similarity between two node embeddings. As shown above, as GNNs inherently evolve the whole graph structure for convolution, it is notoriously challenging for scaling it up. It is worth noting that evolving sufficient neighbors during training is crucial for GNNs. Many studies (Duan et al., 2022; Zou et al., 2019) have shown that full-batch training generally outperforms minibatch for GNNs on multi graph benchmarks. In practice, the lower borderline of batch size for training GNNs is usually thousands. However, when applying it to LMs, it makes the GNN-LM end-toend training intractable, as a text occupies far more GPU memories than an embedding.

Text Embeddings and Language Models. Transforming text in low-dimensional dense embeddings serves as the upstream of textual graph representation learning and has been widely explored in the literature. To generate sentence embeddings with LMs, two commonly-used methods are (i) average pooling (Reimers and Gurevych, 2019) by taking the average of all word embeddings along with attention mask and (ii) taking the embedding of the [CLS] token (Devlin et al., 2018). With the development of pre-trained language models (Devlin et al., 2018; Liu et al., 2019), particular language models (Li et al., 2020; Reimers and Gurevych, 2019) for sentence embeddings have been proposed and shown promising results in various benchmarks (Muennighoff et al., 2022).



Figure 1: The overview of SimTeG. In *stage 1*, we train a LM with lora (Hu et al., 2022) and then generate the embeddings X as the representation of text. In *stage 2*, we train a GNN on top of the embeddings X, along with the graph structure. The two stages are guided with consistent loss function, e.g., link prediction or node classification.

#### 4 SimTeG: Methodology

287

288

290

291

293

296

301

305

307

We propose an extremely simple two-stage training manner that decouples the training of  $gnn(\cdot)$  and  $lm(\cdot)$ . We first finetune lm on  $\mathcal{T}$  with the downstream task loss:

$$Loss_{cls} = \mathcal{L}_{\theta} \big( \phi(lm(\mathcal{T})), \mathbf{Y} \big),$$
$$Loss_{link} = \mathcal{L}_{\theta} \big( \phi(lm(\mathcal{T}_{src}), lm(\mathcal{T}_{dst}) \big), \mathbf{Y} \big), \quad (1)$$

where  $\phi(\cdot)$  is the classifier (left for *node classification*) or similarity function (right for *link prediction*) and Y is the label. After finetuning, we generate node representations X with the finetuned LM *lm*. In practice, we follow Reimers and Gurevych (2019) to perform mean pooling over the output of the last layer of the LM and empirically find that such a strategy is more stable and converges faster than solely taking the <CLS> token embedding as representation (Zhao et al., 2022). In the second stage, we train gnn on (A, X) with the same task. The corresponding loss is computed by replacing  $lm(\mathcal{T})$  with  $gnn(\mathbf{A}, \mathbf{X})$ . The two stage is fully decoupled and one can take advantage of any existing GNN and LM models. We illustrate the two stages in Fig. 1.

310Regularization with PEFT. When fully finetuning311a LM, the inferred features are prone to overfit the312training labels, which results in collapsed feature313space and thus hindering the generalization in GNN314training. Though PEFT was proposed to accelerate315the finetuning process without loss of performance,316in our two-stage finetuning stage, we empirically

find PEFT (Hu et al., 2022; Houlsby et al., 2019; He et al., 2022) could alleviate the overfitting issue to a large extent and thus provide well-regularized node features. See Sec. 5.2 for empirical analysis. In this work, We take the popular PEFT method, lora (Hu et al., 2022), as the instantiation. 317

318

319

320

321

322

323

324

325

326

327

328

329

331

332

333

334

336

337

338

341

342

343

344

345

Feature space without graph structures is already easily differentiable. We plot the two-dimensional feature space computed by T-SNE (Van der Maaten and Hinton, 2008) of X-SimTeG, X-Fix (features generated by pretrained LM without finetuning), and X-OGB regarding labels on OGBN-Arxiv and OGBN-Products in Fig. 2. In detail, we randomly select 100 nodes each with various labels and use T-SNE to compute its two-dimensional features. As shown below, X-SimTeG has a significantly more distinguishable feature space as it captures more semantic information and is finetuned on the downstream dataset. Besides, we find that X-Fix is more distinguishable than X-OGB, which illustrates the inner semantic capture ability of LMs. Furthermore, in comparison with OGBN-Arixv, features in OGBN-Products is visually indifferentiable, indicating the weaker correlation between semantic information and task-specific labels. It accounts for the less improvement of SimTeG on OGBN-Products in Sec. 5.1.

#### **5** Experiments

In the experiments, we aim at answering three research questions as proposed in the introduction



Figure 2: The two-dimensional feature space of X-SimTeG, X-Fix, and X-OGB for OGBN-Arixv, and OGBN-Products. X-SimTeG denotes the features generated by the finetuned LM. different values and shapes refer to different labels on the specific dataset. The feature values are computed by T-SNE. The LM backbone is e5-large (Wang et al., 2022).

(Sec. 1). For a clear statement, we split and reformat them into the following research questions.

**Q1:** How much could SimTeG generally improve the learning of GNNs on node classification and link prediction? **Q2:** Is PEFT a necessity for LM finetuning stage? **Q3:** How sensitive is GNN training to the selection of LMs?

353

354

Datasets. Focusing on two fundamental tasks node classification and link prediction, we con-356 duct experiments on three prestigious benchmarks: 357 OGBN-Arxiv (Arxiv), OGBN-Products (Products), and OGBL-Citation2 (Hu et al., 2020). The former two are for node classification while the latter one is for link prediction. For the former two, we follow the public split, and all text resources are provided by the officials. For the latter one, OGBL-Citation2, as no official text resources are 364 provided, we take the intersection of it and another dataset ogbn-papers100M w.r.t. unified paper ids, which results in a subset of OGBL-Citation2 with about 2.7M nodes. The public split is further updated according to this subset. In comparison, the original OGBL-Citation2 has about 2.9M nodes, which is on par with the TG version, as the public valid and test split occupies solely 2% overall. As a result, we expect roughly consistent 373 performance for methods on the TG version of 374 OGBL-Citation2 and the original one. We introduce the statistics of the three datasets in Table. A9 376

and the details in Appendix A2.1.

**Baselines.** We compare SimTeG with the official features X-OGB (Hu et al., 2020), which is the mean of word embeddings generated by skip-gram (Mikolov et al., 2013). In addition, for node classification, we include another two SOTA methods: X-GIANT (Chien et al., 2021) and GLEM (Zhao et al., 2022). Particularly, X-\* are methods are different at learning node embeddings and any GNN model could be applied in the downstream task for a fair comparison. To make things consistent, we denote our method as X-SimTeG without further specification.

377

378

379

380

381

382

384

386

388

390

391

392

393

394

395

396

397

398

400

401

402

**GNN Backbones.** Aiming at investigating the general improvement of SimTeG, for each dataset, we select two commonly-used baselines GraphSAGE and MLP besides one corresponding SOTA GNN models based on the official leaderboard<sup>1</sup>. For OGBN-Arxiv, we select RevGAT (Li et al., 2021); for OGBN-Products, we select SAGN+SCR (Sun et al., 2021; Zhang et al., 2021a); and for ogbn-citation2, we select SEAL (Zhang and Chen, 2018).

**LM Backbones.** For retrieval LM backbones, we select three popular LMs on MTEB (Muennighoff et al., 2022) leaderboard<sup>2</sup> w.r.t. model size and

<sup>&</sup>lt;sup>1</sup>https://ogb.stanford.edu/docs/leader\_ nodeprop

<sup>&</sup>lt;sup>2</sup>https://huggingface.co/spaces/mteb/

Table 1: The performance of SOTA GNN, GraphSAGE and MLP on OGBN-Arxiv, OGBN-Products, OGBL-Citation2, which are averaged over 10 runs (Please note the we solely train LM once to generate the node embeddings). The results of GLEM is from the orignal paper. We **bold** the best results w.r.t. the same GNN backbone and red color the smallest  $\Delta_{MLP}$  and  $\Delta_{GNN}$ .

Dataset	Metric	Method	SOTA GNN	A 2-la	yer Simp	ole MLP / GNN	
Dutuset	1010tile	1010tilot	RevGAT	MLP	$\Delta_{MLP}$	GraphSAGE	$\Delta_{GNN}$
		X-OGB	$74.01 \pm 0.29$	$47.73\pm0.29$	25.24	$71.80\pm0.20$	3.40
Arviv	$\Lambda cc$ (%)	X-GIANT	$75.93 \pm 0.22$	$71.08\pm0.22$	4.85	$73.70\pm0.09$	2.23
AIXIV	Acc. (70)	GLEM	$76.97\pm0.19$	-	-	$75.50 \pm 0.24$	1.47
		X-SimTeG	$ $ 77.04 $\pm$ 0.13	$\begin{array}{ }\textbf{74.06}\pm\textbf{0.13}\end{array}$	2.98	$\textbf{76.84} \pm \textbf{0.34}$	0.20
Dataset	Metric	Metric Method		A 2-la	yer Simp	le MLP / GNN	
Dutuset	Dutaset metric		SAGN+SCR	MLP	$\Delta_{MLP}$	GraphSAGE	$\Delta_{GNN}$
		X-OGB	$81.82 \pm 0.44$	$50.86 \pm 0.26$	30.96	$78.81 \pm 0.23$	3.01
Producto	$\Lambda_{00}$ ( $\mathcal{O}_{2}$ )	X-GIANT	$86.12 \pm 0.34$	$\textbf{77.58} \pm \textbf{0.24}$	8.54	$82.84\pm0.29$	3.28
Tiouucis	Acc. (70)	GLEM	$\textbf{87.36} \pm \textbf{0.07}$	-	-	$83.16\pm0.19$	4.20
		X-SimTeG	$85.40 \pm 0.28$	$76.73 \pm 0.44$	8.67	$\textbf{84.59} \pm \textbf{0.44}$	0.81
Dataset	Metric	Method	SOTA GNN	A 2-la	yer Simp	le MLP / GNN	
	Dataset metric		SEAL	MLP	$\Delta_{MLP}$	GraphSAGE	$\Delta_{GNN}$
		X-OGB	$86.14 \pm 0.40$	$25.44\pm0.01$	60.70	$77.31 \pm 0.90$	8.83
Citation2	<b>MKK</b> (%)	X-SimTeG	$\textbf{86.66} \pm \textbf{1.21}$	$\textbf{72.90} \pm \textbf{0.14}$	13.76	$\textbf{85.13} \pm \textbf{0.73}$	1.53
Chatlon2	$\mathbf{H}_{ta} \otimes 2 \left( \mathcal{O} \right)$	X-OGB	$90.92 \pm 0.32$	$28.22\pm0.02$	62.70	$85.56 \pm 0.69$	5.36
	nits@3(%)	X-SimTeG	$\textbf{91.42} \pm \textbf{0.19}$	$\textbf{80.55} \pm \textbf{0.13}$	10.87	$\textbf{91.62} \pm \textbf{0.87}$	-0.20

403

404

405

412

413 414

415

416

417

418

419

420

421

422

423

424

performance on classification and retrieval: all-MiniLM-L6-v2 (Reimers and Gurevych, 2019), allroberta-large-v1 (Reimers and Gurevych, 2019), and e5-large-v1 (Wang et al., 2022). We present the properties of the three LMs in Table. A11. Hyperparameter search. We utilize optuna (Akiba et al., 2019) to perform hyperparameter search on all tasks. The search space for LMs and GNNs on all datasets is presented in Appendix A2.3.

#### 5.1 **Q1:** How much could SimTeG generally improve the learning of GNNs on node classification and link prediction?

In this section, we conduct experiments to show the superiority of SimTeG on improving the learning of GNNs on node classification and link prediction. The reported results are selected based on the validation dataset. We present the results based on e5-large backbone in Table. 1 and present the comprehensive results of node classification and link prediction with all the three selected backbones in Table A5 and Table A6. Specifically, in Table 1, we present two comparison metric  $\Delta_{MLP}$ 

leaderboard

and  $\Delta_{GNN}$  to describe the performance margin of (SOTA GNN, MLP) (SOTA GNN, GraphSAGE), respectively. The smaller the value is, even negative, the better the performance of simple models is. In addition, we ensemble the GNNs with multiple node embeddings generated by various LMs and text resources on OGBN-Arxiv and show the results in Table 2. We find several interesting observations as follows.

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

**Observation 1:** (*X-SimTeG* + *GraphSAGE*) consistently outperforms (X - OGB + SOTA GNN) on all the three datasets. This finding implies that the incorporation of advanced text features can bypass the necessity of complex GNNs, which is why we perceive our method to be *frustratingly* simple. Furthermore, when replacing GraphSAGE with the corresponding SOTA GNN in X-SimTeG, although the performance is improved moderately, this margin of improvement is notably smaller compared to the performance gap on X-OGB. Particularly, we show that the simple 2-layer GraphSAGE achieves comparable performance with the dataset-specific SOTA GNNs. Particularly, on OGBN-Arxiv, Graph-SAGE achieves 76.84%, taking the *4-th* place in the corresponding leaderboard (by 2023-08-01).

Table 2: The performance of GraphSAGE and RevGAT trained on OGBN-Arxiv with ad	ditional text	t attributes
provided by He et al. (2023). LMs for ensembling are e5-large and all-roberta-large-v1. We	select the to	pp-3 SOTA
methods from the leaderboard of OGBN-Arxiv (accessed on 2023-07-18) for comparison and	gray color	our results
(reported over 10 runs).		

Rank	Method	GNN Backbone	Valid Acc. (%)	Test Acc. (%)
1	TAPE + SimTeG (ours)	RevGAT	$\textbf{78.46} \pm \textbf{0.04}$	$\textbf{78.03} \pm \textbf{0.07}$
2	<b>TAPE</b> (He et al., 2023)	RevGAT	$77.85\pm0.16$	$77.50\pm0.12$
3	TAPE + SimTeG (Ours)	GraphSAGE	$77.89 \pm 0.08$	$77.48 \pm 0.11$
4	GraDBERT (Mavromatis et al., 2023)	RevGAT	$77.57\pm0.09$	$77.21\pm0.31$
5	GLEM (Zhao et al., 2022)	RevGAT	$77.46 \pm 0.18$	$76.94\pm0.25$

Besides, on OGBL-Citation2, GraphSAGE even 450 outperforms the SOTA method SEAL on Hits@3. 451 **Observation 2: With additional text attributes,** 452 SimTeG with Ensembling achieves new SOTA 453 performance on OGBN-Arxiv. We further demon-454 strate the effectiveness of SimTeG by ensembling 455 the node embeddings generated by different LMs 456 and texts. For text, we use both the original text 457 provided by Hu et al. (2020) and the additional 458 text attributes<sup>3</sup> provided by He et al. (2023), which 459 is generated by ChatGPT. For LMs, we use both 460 e5-large and all-roberta-large-v1. We train Graph-461 SAGE or RevGAT on those node embeddings gen-462 463 erated by various LMs and texts, and make the final predictions with weighted ensembling (taking the 464 weighted average of all predictions). As shown 465 in Table 2, with RevGAT, we achieve new SOTA 466 performance on OGBN-Arxiv with 78.03% test ac-467 curacy, more than 0.5 % higher than the previous 468 SOTA performance (77.50%) achieved by He et al. 469 (2023). It further validates the importance of text 470 features and the effectiveness of SimTeG. 471

**Observation 3: Text attributes are unequally** 472 important for different datasets. As shown in Ta-473 ble 1, we compute  $\Delta_{MLP}$  which is the performance 474 gap between MLP and SOTA GNNs. Empirically, 475 this value indicates the importance of text attributes 476 on the corresponding dataset, as MLP is solely 477 trained on the texts (integrated with SOTA LMs) 478 while SOTA GNN additionally takes advantage of 479 graph structures. Therefore, approximately, the 480 less  $\Delta_{MLP}$  is, the more important text attributes are. 481 As presented in Table 1,  $\Delta_{MLP}$  on OGBN-Arxiv 482 is solely 2.98, indicating the text attributes are 483 more important, in comparison with the ones in 484 OGBN-Products and OGBL-Citation2. This em-485

pirically indicates why the performance of SimTeG in OGBN-Products does not perform as well as the one in OGBN-Arxiv. We show a sample of text in OGBN-Arxiv and OGBN-Products respectively in Appendix A2.1. We find that the text in OGBN-products resembles more a bag of words, which account for the less improvement when using LM features.

# 5.2 Q2: Is PEFT a necessity for LM finetuning stage?

In this ablation study, we analyze the effectiveness of PEFT for LM finetuning stage in SimTeG. We summarize the training, validation, and test accuracy of two stages: LM finetuning stage and GNN training stage. The results of node classification are presented in Table 4.

Observation 4: PEFT could significantly alleviate the overfitting problem during finetuning LM and further facilitate the training of GNNs with regularized features. As shown in Table 4, due to the excessively strong learning capacity of LMs, finetuning LMs on the downstream task causes a severe overfitting problem. Although fullfinetuning outperforms PEFT in LM stage, training GNNs on the derived features gains notably less improvement. In contrast, PEFT could significantly mitigate the overfitting issue according to  $\Delta_{overfit}$ in LM finetuning stage and assist the training of GNNs with regularized features to gain considerable improvement compared with full-finetuning.

# 5.3 Q3: How sensitive is GNN training to the selection of LMs?

**Observation 5: GNN's training is moderately sensitive to the selection of LMs.** We select three retrieval LMs based on their rank in MTEB leaderboard in terms of the classification and retrieval performance. Interestingly, based on the leaderboard, the performance ranking is *e5-large* > *all-roberta*-

 $<sup>^{3}</sup>$ It is worth noting that as GPT-4 used by He et al. (2023) does not release their training recipe, we do not know whether the arxiv papers are included during training, which may lead to a label leakage problem.



Figure 3: (*Left*): The performance of GraphSAGE trained on SimTeG with different LM backbones. (*Right*): GNN's performance on OGBN-Arxiv using LMs of various sizes, indicated by the bubble size.

Table 3: The performance of Graph and MLP trained on SimTeG backed with all-roberta-large-v1 and roberta-large, which have the same model architecture. we **bold** the best results for each comparison in X-Fix and X-SimTeG. all results are reported based on 10 runs.

datasets	Metric	X_type		<b>X</b> -Fix	X-SimTeG		
		LM Backbone	roberta-large	all-roberta-large-v1	roberta-large	all-roberta-large-v1	
Arxiv	Acc.	MLP GraphSAGE	$ \begin{vmatrix} 61.15 \pm 0.83 \\ 72.15 \pm 0.59 \end{vmatrix} $	$\begin{array}{c} \textbf{72.58} \pm \textbf{0.25} \\ \textbf{75.51} \pm \textbf{0.23} \end{array}$	$ \begin{vmatrix} 71.55 \pm 0.24 \\ 75.48 \pm 0.16 \end{vmatrix} $	$\begin{array}{c} \textbf{74.32} \pm \textbf{0.12} \\ \textbf{76.18} \pm \textbf{0.37} \end{array}$	
Products	Acc.	MLP GraphSAGE	$ \begin{vmatrix} 68.14 \pm 0.23 \\ 77.65 \pm 0.34 \end{vmatrix} $	$\begin{array}{c} \textbf{70.10} \pm \textbf{0.08} \\ \textbf{82.38} \pm \textbf{0.60} \end{array}$	$ \begin{vmatrix} 78.45 \pm 0.14 \\ 83.56 \pm 0.21 \end{vmatrix} $	$\begin{array}{c} \textbf{77.48} \pm \textbf{0.19} \\ \textbf{83.68} \pm \textbf{0.32} \end{array}$	
Citation2	MRR	MLP GraphSAGE	$ \begin{vmatrix} 00.20 \pm 0.01 \\ 79.71 \pm 0.27 \end{vmatrix} $	$\begin{array}{c} \textbf{70.12} \pm \textbf{0.12} \\ \textbf{83.20} \pm \textbf{0.40} \end{array}$	$\begin{vmatrix} 63.15 \pm 0.20 \\ 84.37 \pm 0.34 \end{vmatrix}$	$\begin{array}{c} \textbf{72.90} \pm \textbf{0.14} \\ \textbf{85.13} \pm \textbf{0.73} \end{array}$	

Table 4: The training results of finetuning LM (*LM* stage) and further training GNN on top of derived features (*GNN stage*). We compare the results of PEFT (SimTeG) with full-finetuning (*X*-FULL). The LM backbone is e5-large and the GNN backbone is Graph-SAGE. We **bold** the better results on each comparison.  $\Delta_{overfit}$  computes (Train Acc. - Test Acc.) to measure the overfitting.

datasets	Stage	X_type	Train Acc.	Valid Acc	Test Acc.	$\Delta_{overfit}$
	тм	X-FULL	82.33	75.85	74.77	7.56
Arxiv	LIVI	X-SimTeG	75.72	75.40	74.31	1.41
	GNN	X-FULL	84.39	76.73	75.28	9.11
	UNIN	X-SimTeG	79.37	77.47	76.85	2.52
	тм	X-FULL	95.46	91.70	78.70	16.76
Products	LIVI	X-SimTeG	89.45	88.85	77.81	11.64
	GNN	X-FULL	96.42	93.18	81.80	14.62
	UNIN	X-SimTeG	95.37	93.57	84.58	10.79

*large-v1* > *all-MiniLM-L6-v2*, which is consistent with their overall performance in left subfigure Figure 3. Based on the right subfigure of Figure 3, we find that the performance of downstream GNN is not closely correlated with the LM size, but probably with ability to generate representative text embeddings. To further validate this, we perform an ablation study regarding the comparison between a pretrained LM and the same LM finetuned for retrieval tasks. The results are shown in Table 3. We observe that given the same architecture, the models specifically finetuned for retrieval tasks (*all-roberta-large-v1*) generally perform better on tasks of TG representation learning.

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

#### 6 Conclusion

In this work, we propose a frustratingly simple approach SimTeG for TG representation learning. We show that with a parameter-efficiently finetuned LM on the same downstream task first, a simple two-layer GraphSAGE trained on the generated node embeddings can achieve on-par state-of-theart (SOTA) performance on OGBN-Arxiv (77.48 %). It indicates that when incorporating advaced text features, one can bypass the necessity of using complex GNN architectures and the combination of LM + Simple GNN is capable of achieving satisfactory results on graph tasks including node classification and link prediction.

531

#### References

552

557

558

559

565

568

571

572

573

574

575

577

579

580

583

584

585

586

589

590

591

593

599

601

604

- Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. 2019. Optuna: A nextgeneration hyperparameter optimization framework. In Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining, pages 2623–2631.
- Dexiong Chen, Leslie O'Bray, and Karsten Borgwardt. 2022. Structure-aware transformer for graph representation learning. In *International Conference on Machine Learning*, pages 3469–3489. PMLR.
- Eli Chien, Wei-Cheng Chang, Cho-Jui Hsieh, Hsiang-Fu Yu, Jiong Zhang, Olgica Milenkovic, and Inderjit S Dhillon. 2021. Node feature extraction by self-supervised multi-scale neighborhood prediction. *arXiv preprint arXiv:2111.00064*.
  - Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Keyu Duan, Zirui Liu, Peihao Wang, Wenqing Zheng, Kaixiong Zhou, Tianlong Chen, Xia Hu, and Zhangyang Wang. 2022. A comprehensive study on large-scale graph training: Benchmarking and rethinking. In *Advances in Neural Information Processing Systems*, volume 35, pages 5376–5389.
- Will Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. In *NeuIPS*, pages 1024–1034.
- Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. 2022. Towards a unified view of parameter-efficient transfer learning. In International Conference on Learning Representations.
- Xiaoxin He, Xavier Bresson, Thomas Laurent, and Bryan Hooi. 2023. Explanations as features: Llmbased features for text-attributed graphs. *arXiv preprint arXiv:2305.19523*.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019.
   Parameter-efficient transfer learning for nlp. In *International Conference on Machine Learning*, pages 2790–2799. PMLR.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference* on Learning Representations.
- Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. 2020. Open graph benchmark: Datasets for machine learning on graphs. Advances in neural information processing systems, 33:22118–22133.

Lianzhe Huang, Dehong Ma, Sujian Li, Xiaodong Zhang, and Houfeng Wang. 2019. Text level graph neural network for text classification. *arXiv preprint arXiv:1910.02356*.

606

607

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

- Md Shamim Hussain, Mohammed J Zaki, and Dharmashankar Subramanian. 2022. Global self-attention as a replacement for graph convolution. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 655–665.
- Bowen Jin, Wentao Zhang, Yu Zhang, Yu Meng, Xinyang Zhang, Qi Zhu, and Jiawei Han. 2023. Patton: Language model pretraining on text-rich networks. *arXiv preprint arXiv:2305.12268*.
- Thomas N Kipf and Max Welling. 2016. Semisupervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*.
- Bohan Li, Hao Zhou, Junxian He, Mingxuan Wang, Yiming Yang, and Lei Li. 2020. On the sentence embeddings from pre-trained language models. *arXiv preprint arXiv:2011.05864*.
- Guohao Li, Matthias Müller, Bernard Ghanem, and Vladlen Koltun. 2021. Training graph neural networks with 1000 layers. In *International conference on machine learning*, pages 6437–6449. PMLR.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692.
- Costas Mavromatis, Vassilis N Ioannidis, Shen Wang, Da Zheng, Soji Adeshina, Jun Ma, Han Zhao, Christos Faloutsos, and George Karypis. 2023. Train your own gnn teacher: Graph-aware distillation on textual graphs. *arXiv preprint arXiv:2304.10668*.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems*, 26.
- Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. 2022. Mteb: Massive text embedding benchmark. *arXiv preprint arXiv:2210.07316*.
- OpenAI. 2023. Introducing chatgpt. https:// openai.com/blog/chatgpt. Accessed: 2023-07-18.
- Wonpyo Park, Woong-Gi Chang, Donggeon Lee, Juntae Kim, et al. 2022. Grpe: Relative positional encoding for graph transformer. In *ICLR2022 Machine Learning for Drug Discovery*.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*.

- Chuxiong Sun, Hongming Gu, and Jie Hu. 2021. Scalable and adaptive graph neural networks with self-label-enhanced training. *arXiv preprint arXiv:2104.09376*.
  - Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-sne. *Journal of machine learning research*, 9(11).
  - Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2022. Text embeddings by weakly-supervised contrastive pre-training. *arXiv preprint arXiv:2212.03533*.
  - Luyu Wang, Yujia Li, Ozlem Aslan, and Oriol Vinyals. 2021. Wikigraphs: A wikipedia textknowledge graph paired dataset. *arXiv preprint arXiv:2107.09556*.

671

672

673

674

675

676

677

678

679

688

690

695

704

705 706

707

711 712

- Zhanghao Wu, Paras Jain, Matthew Wright, Azalia Mirhoseini, Joseph E Gonzalez, and Ion Stoica. 2021.
   Representing long-range context for graph neural networks with global attention. Advances in Neural Information Processing Systems, 34:13266–13279.
- Junhan Yang, Zheng Liu, Shitao Xiao, Chaozhuo Li, Defu Lian, Sanjay Agrawal, Amit Singh, Guangzhong Sun, and Xing Xie. 2021. Graphformers: Gnn-nested transformers for representation learning on textual graph. *Advances in Neural Information Processing Systems*, 34:28798–28810.
- Zhilin Yang, William Cohen, and Ruslan Salakhudinov. 2016. Revisiting semi-supervised learning with graph embeddings. In *International conference on machine learning*, pages 40–48. PMLR.
- Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin Zheng, Guolin Ke, Di He, Yanming Shen, and Tie-Yan Liu. 2021. Do transformers really perform badly for graph representation? *Advances in Neural Information Processing Systems*, 34:28877–28888.
- Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L Hamilton, and Jure Leskovec. 2018.
  Graph convolutional neural networks for web-scale recommender systems. In *Proceedings of the 24th* ACM SIGKDD international conference on knowledge discovery & data mining, pages 974–983.
- Hanqing Zeng, Hongkuan Zhou, Ajitesh Srivastava, Rajgopal Kannan, and Viktor Prasanna. 2019. Graphsaint: Graph sampling based inductive learning method. arXiv preprint arXiv:1907.04931.
- Chenhui Zhang, Yufei He, Yukuo Cen, Zhenyu Hou, and Jie Tang. 2021a. Improving the training of graph neural networks with consistency regularization. *arXiv preprint arXiv:2112.04319*.
- Jiong Zhang, Wei-Cheng Chang, Hsiang-Fu Yu, and Inderjit Dhillon. 2021b. Fast multi-resolution transformer fine-tuning for extreme multi-label text classification. Advances in Neural Information Processing Systems, 34:7267–7280.

Muhan Zhang and Yixin Chen. 2018. Link prediction based on graph neural networks. In Advances in Neural Information Processing Systems, pages 5165– 5175. 713

714

715

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

- Yufeng Zhang, Xueli Yu, Zeyu Cui, Shu Wu, Zhongzhen Wen, and Liang Wang. 2020. Every document owns its structure: Inductive text classification via graph neural networks. *arXiv preprint arXiv:2004.13826*.
- Jianan Zhao, Meng Qu, Chaozhuo Li, Hao Yan, Qian Liu, Rui Li, Xing Xie, and Jian Tang. 2022. Learning on Large-scale Text-attributed Graphs via Variational Inference. *arXiv preprint arXiv:2210.14709*.
- Jason Zhu, Yanling Cui, Yuming Liu, Hao Sun, Xue Li, Markus Pelger, Tianqi Yang, Liangjie Zhang, Ruofei Zhang, and Huasha Zhao. 2021. Textgnn: Improving text encoder via graph neural network in sponsored search. In *Proceedings of the Web Conference 2021*, pages 2848–2857.
- Difan Zou, Ziniu Hu, Yewen Wang, Song Jiang, Yizhou Sun, and Quanquan Gu. 2019. Layer-dependent importance sampling for training deep and large graph convolutional networks. *Advances in neural information processing systems*, 32.

### A1 More Experiment Results

#### A1.1 Comprehensive Results of Main Experiments

Table A5: Node Classification Accuracy of X-SimTeG on ogbn-arxiv (Arxiv) and ogbn-products (Products). All reported results are averaged over 10 runs in the format of mean  $\pm$  std. We red color the best results and blue color the runner-ups with the same GNN backbone.  $\uparrow$  (%) denotes the improvement of X-SimTeG over the original feature X-OGB.  $\Delta_{MLP}$  and  $\Delta_{GNN}$  denotes the extreme value difference among all methods (including MLP) and GNNs, respectively.

		1	1								
Datasets	GNN	Acc. (%)		Baselines				X-SimTe	eG		
			X-OGB	X-GIANT	$GLEM^a$	MiniLM-L6	$\uparrow$ (%)	e5-large	$\uparrow$ (%)	roberta-large	$\uparrow$ (%)
	MUD	val	$49.14 \pm 0.27$	$72.02\pm0.16$	-	$71.59 \pm 0.07$	22.45	$75.08 \pm 0.09$	26.66	$74.80\pm0.07$	25.66
	MLP	test	$47.73 \pm 0.29$	$71.08\pm0.22$	-	$70.56 \pm 0.09$	22.83	$74.06\pm0.13$	26.33	$74.32\pm0.12$	26.59
	CurrleACE	val	$72.80 \pm 0.18$	$74.58\pm0.20$	$76.45\pm0.05$	$75.92 \pm 0.17$	3.12	$77.47\pm0.14$	4.67	$76.86\pm0.13$	4.06
	GraphSAGE	test	$71.80\pm0.20$	$73.70\pm0.09$	$75.50\pm0.24$	$75.14 \pm 0.30$	3.34	$76.84 \pm 0.34$	5.04	$76.18\pm0.37$	4.38
Arxiv	CAMER	val	$71.49\pm0.41$	$76.36\pm0.09$	$76.95\pm0.14$	$76.75 \pm 0.11$	5.26	$77.90 \pm 0.12$	6.41	$77.57\pm0.15$	6.08
	GAMLF	test	$70.61 \pm 0.52$	$75.26\pm0.15$	$75.62\pm0.23$	$75.46 \pm 0.17$	4.85	$76.92\pm0.10$	6.31	$76.72\pm0.19$	6.11
	SACN	val	$72.74 \pm 0.39$	$75.76\pm0.21$	-	$76.84 \pm 0.08$	4.10	$78.03\pm0.05$	5.29	$77.63\pm0.16$	4.89
	SAGN	test	$71.76 \pm 0.41$	$74.39\pm0.38$	-	$75.50 \pm 0.23$	3.74	$76.85\pm0.12$	5.09	$76.59\pm0.17$	4.83
	DerecAT	val	$75.10 \pm 0.15$	$76.97\pm0.08$	$77.49 \pm 0.17$	$76.86 \pm 0.24$	1.76	$77.68 \pm 0.07$	2.58	$76.32\pm0.18$	1.22
	RevGAI	test	$74.01 \pm 0.29$	$75.93\pm0.22$	$76.97\pm0.19$	$75.96\pm0.21$	1.95	$77.04\pm0.13$	3.03	$75.88\pm0.58$	1.87
	$\Delta_{MLP}/2$	$\Delta_{GNN}$	25.24 / 3.40	4.85 / 2.23	-	5.40 / 0.82	-	2.98 / 0.20	-	2.40 / 0.84	-
	MUD	val	$63.44 \pm 0.30$	$89.67\pm0.07$	-	$86.82 \pm 0.02$	23.38	$88.75 \pm 0.04$	25.31	$90.01\pm0.03$	26.57
	MLP	test	$50.86 \pm 0.26$	$77.58 \pm 0.24$	-	$72.36 \pm 0.12$	21.50	$76.73 \pm 0.44$	25.87	$77.48 \pm 0.19$	26.62
	Carabeace	val	$90.03 \pm 0.08$	$93.49\pm0.09$	$93.84 \pm 0.12$	$93.49 \pm 0.08$	3.46	$93.57 \pm 0.20$	3.54	$93.34 \pm 0.09$	3.31
Decducto	GraphSAGE	test	$78.81 \pm 0.23$	$82.84\pm0.29$	$83.16\pm0.19$	$82.04 \pm 0.57$	3.23	$84.59\pm0.44$	5.78	$83.68\pm0.32$	4.87
FIODUCIS	SACNISCE	val	$91.83 \pm 0.24$	$94.04\pm0.12$	$94.00\pm0.03$	$92.89 \pm 0.07$	1.06	$94.12\pm0.10$	2.29	$94.13\pm0.12$	2.30
	SAGIN+SUK	test	$81.82 \pm 0.44$	$86.12\pm0.34$	$87.36\pm0.07$	$82.43 \pm 0.40$	0.61	$85.40 \pm 0.28$	3.58	$85.23\pm0.32$	3.41
	$\Delta_{MLP}/2$	GNN	30.96 / 3.01	8.54/3.28	-	10.07 / 0.39		8.67 / 0.81	-	7.75 / 1.55	-

<sup>a</sup> results are from the original papers.

#### A1.2 P value analysis

As shown in the upper sub-table in Table A7, all p values are lower than 0.05, indicating the significant improvement of SimTeG. It is worth noting that as the results of GLEM are reported by the original paper and we do not have the results for each individual experiment, we are not able to compute the corresponding p values. We do acknowledge that there is a subtle difference between SimTEG and GLEM and GLEM outperforms SimTEG on OGBN-Products. This phenomenon is discussed in our Observation 4 in Section 5.1 of the paper.

In addition, as shown in the bottom sub-table in Table A7, the p values of SimTeG are significantly smaller than the baseline embeddings. Specifically, the p values of SimTeG on OGBN-Arxiv and OGBL-Citation2 are close or larger than 0.05. This further supports our key findings: in cooperation with advanced text embeddings, one can bypass the necessity of using complex GNN models.

#### A1.3 Comparison with More LM-involved Methods

The results of GraphFormer on OGBN-Arxiv and OGBN-Products are directly borrowed from the GLEM paper (Zhao et al., 2022). since the datasets and split are exactly the same. We run GraphFormer on ogbl-citation2 for 10 times and report the mean with std. For the hyperparameter setting, we use the default parameters, and the batch size is set to 100 to make it consistent with the reported results in GLEM. As shown in the table, SimTeG performs consistently better than GraphFormer. It is possibly because (i) the GNN-nested architecture of GraphFormer solely allows 1-hop message passing, which limits the express ability of GNN models; (ii) GraphFormer's implementation modifies the architecture code of BERT (Devlin et al., 2018) and cannot be easily extended to other SOTA embedding models nowadays.

#### A1.4 More Results of Ablation Studies

Towards a comprehensive understanding of the effectiveness of SimTeG, we further investigate the convergence of GNNs with SimTeG. We compare the training convergence and the corresponding validation performance of GNNs trained on SimTeG, *X*-OGB, and *X*-FIX, where *X*-FIX denotes the node embeddings generated by the pretrained LMs without finetuning. The illustration is placed in Fig. A4.

Table A6: Link prediction results on *OGBL-Citation2-2.7M* (Citation2). All reported results are averaged over 10 runs. We red color the best results and blue color the runner-ups with the same GNN backbone.  $\uparrow$  (%) denotes the improvement of **X**-SimTeG over the original feature **X**-OGB.  $\Delta_{MLP}$  and  $\Delta_{GNN}$  denotes the margin of (MLP, SEAL) and (GraphSAGE, SEAL), respectively. We use blue color denoting the negative values and red denoting positive. Specifically, in the context of  $\Delta$ , positives indicate MLP/GraphSAGE performs better than SEAL.

Metrics	GNN	Split	Baselines			X-SimTe	eG		
metries	GINI	opin	X-OGB	MiniLM-L6	$\uparrow$ (%)	roberta-large	$\uparrow$ (%)	e5-large	$\uparrow$ (%)
	МГР	val	$25.37\pm0.09$	$64.56\pm0.15$	39.19	$70.20 \pm 0.19$	44.83	$72.79 \pm 0.17$	47.42
	MLP	test	$25.44 \pm 0.01$	$64.49 \pm 0.18$	39.05	$70.32\pm0.22$	44.88	$72.90\pm0.14$	47.46
MDD	GraphsACE	val	$77.40 \pm 0.88$	$83.13\pm0.72$	5.73	$85.27\pm0.78$	7.87	$85.20\pm0.69$	7.80
MKK	GraphSAGE	test	$77.31\pm0.90$	$83.09\pm0.75$	5.78	$85.29\pm0.70$	7.98	$85.13\pm0.73$	7.82
	SEAL	val	$87.21\pm0.03$	$88.33\pm0.30$	1.12	$88.29 \pm 0.45$	1.08	$88.56 \pm 0.38$	1.35
	SEAL	test	$86.14\pm0.40$	$86.69\pm0.43$	0.55	$87.02\pm0.46$	0.88	86.66 ± 1.21	0.52
	$\Delta_{MLP}/\Delta_G$	NN	-60.70/-8.83	-22.20/-3.60	-	-16.70/-1.73	-	-13.76/-1.53	-
	MLD	val	$15.04\pm0.09$	$52.29\pm0.18$	37.25	59.46 ± 0.19	44.42	$62.21 \pm 0.23$	47.17
	MLP	test	$15.11\pm0.06$	$52.18\pm0.25$	37.07	$59.66 \pm 0.26$	44.55	$62.31\pm0.19$	47.20
Uite@1	GraphsACE	val	$67.28 \pm 1.20$	$74.83 \pm 1.02$	7.55	$77.98 \pm 1.20$	10.70	$77.73 \pm 0.89$	10.45
mus@1	GIAPHSAGE	test	$67.09 \pm 1.25$	$74.79 \pm 1.10$	7.70	$77.99 \pm 0.89$	10.90	$77.66 \pm 0.91$	10.57
	SEAL	val	$82.76\pm0.14$	$84.35\pm0.42$	1.59	$84.25\pm0.79$	1.49	$84.70\pm0.58$	1.94
	SEAL	test	$81.74\pm0.46$	$81.40\pm0.96$	-0.34	$82.34\pm0.79$	0.60	$81.15 \pm 2.04$	-0.59
	$\Delta_{MLP}/\Delta_G$	NN	-66.63 / -14.65	-29.22 /-6.61	-	-22.68 / -4.35	-	-18.84/-3.39	-
	MLD	val	$28.06\pm0.10$	$72.60\pm0.16$	44.54	$77.56 \pm 0.23$	49.50	$80.42 \pm 0.15$	52.36
	MLF	test	$28.22\pm0.02$	$72.62\pm0.19$	44.40	$77.66\pm0.24$	49.44	$80.55\pm0.13$	52.33
Hits@3	GraphSAGE	val	$85.54\pm0.69$	$90.17\pm0.61$	4.63	$91.55\pm0.98$	6.01	$91.72 \pm 0.90$	6.18
mis@J	GraphSAGE	test	$85.56\pm0.69$	$90.16\pm0.51$	4.60	$91.57 \pm 1.10$	6.01	$91.62\pm0.87$	6.06
	SEAL	val	$91.36\pm0.44$	$92.00\pm0.07$	0.64	$92.15\pm0.19$	0.79	$91.75\pm0.18$	0.39
	SEAL	test	$90.92\pm0.32$	$91.42\pm0.60$	0.50	$91.52\pm0.56$	0.60	$91.42\pm0.19$	0.50
	$\Delta_{MLP}/\Delta_G$	NN	-62.70/-5.36	-18.80/-1.26	-	-13.86 / 0.05	-	-10.87 / 0.20	-
	MLD	val	$46.73\pm0.14$	$87.62\pm0.06$	40.89	89.80 ± 0.20	43.07	$91.74 \pm 0.08$	45.01
	MLF	test	$46.59\pm0.11$	$87.57\pm0.12$	40.98	$89.66 \pm 0.14$	43.07	$91.74\pm0.10$	45.15
Uite@10	GraphsACE	val	$94.29\pm0.19$	$96.25\pm0.13$	1.96	$96.61\pm0.12$	2.32	$96.71 \pm 0.09$	2.42
mus@10	OTAPIISAGE	test	$94.37\pm0.17$	$96.30\pm0.13$	1.93	$96.64\pm0.12$	2.27	$96.74\pm0.11$	2.37
	SEAL	val	$94.59\pm0.14$	$94.88\pm0.25$	0.29	$95.08 \pm 0.12$	0.49	$95.08 \pm 0.21$	0.49
	JEAL	test	$93.90\pm0.49$	$94.40\pm0.07$	0.50	$93.95 \pm 0.37$	0.05	$94.54\pm0.25$	0.64
	$\Delta_{MLP}/\Delta_G$	NN	-47.31 / -0.47	-6.83 /1.90	-	-4.29 / 2.66	-	-2.80 / 2.20	-

763It is worth noting that we use the training accuracy on OGBN-Arxiv and OGBN-Products to denote their764convergence since we utilize *label smoothing* during training which make the training loss not directly765comparable on them. Based on Fig. A4, we have the following observation:

**Observation 6: SimTeG moderately speeds up and stabilizes the training of GNNs.** As shown in Fig. A4, GNNs with SimTeG generally converge faster than the ones with X-OGB and X-FIX. With SimTeG, GraphSAGE could converge within 2 epochs on OGBN-Arxiv and OGBN-Products. In contrast, training on the features directly generated by the pretrained LMs (i.e., X-FIX) converges much slower, even slower than one of X-OGB (possibly due to a larger hidden dimension). This further indicates the benefits of SimTeG.

## A2 Reproducibility Statement

768

771

772

774

775

## A2.1 Details of TG Version for the three OGB datasets

In this section, we present the details of the TG version of OGBN-Arxiv, OGBN-Products, and OGBL-Citation2. The statistics of the three datasets are shown in Table A9 and the text resources are shown in Table A10.

OGBN-Arxiv. OGBN-Arxiv is a directed academic graph, where node denotes papers and edge denotes directed citation. The task is to predict the category of each paper as listed in https://arxiv.org. For its TG version, we use the same split as Hu et al. (2020). The text for each node is its title and abstract. We concatenate them for each node with the format of "*title: {title}; abstract: {abstract}*" as the corresponding

Table A7: The p values for two comparisons, SimTeG v.s. baseline (GIANT/OGB) and GraphSAGE v.s. SOTA GNN on three datasets. p value smaller than 0.05 means that SimTeG (or SOTA GNN) is significantly better than the baseline (GraphSAGE)

Dataset	GNN	SimTeG	Baseline	<b>P-Value</b>	P <0.05
OCDN Amin	RevGAT	77.04	75.93	7.77e-14	True
OGBN-AIXIV	GraphSAGE	76.84	73.70	4.11e-17	True
OCDN Dro du oto	SAGN+SLE	85.40	86.12	4.15e-06	True
OGBN-Products	GraphSAGE	84.59	82.84	9.01e-10	True
OCDL Citation2	SEAL	91.42	90.92	0.0023	True
OGBL-Citation2	GraphSAGE	91.62	85.13	5.01e-12	True
Dataset	Embeddings	GraphSAGE	SOTA GNN	P-Value	P <0.05
Dataset	Embeddings X-SimTeG	GraphSAGE	<b>SOTA GNN</b> 77.04	<b>P-Value</b> 0.0427	<b>P &lt;0.05</b> True
Dataset OGBN-Arxiv	Embeddings X-SimTeG X-GIANT	<b>GraphSAGE</b> 76.84 73.70	<b>SOTA GNN</b> 77.04 75.93	<b>P-Value</b> 0.0427 7.79e-22	<b>P &lt;0.05</b> True True
Dataset OGBN-Arxiv OGBN Brachusta	Embeddings X-SimTeG X-GIANT X-SimTeG	GraphSAGE 76.84 73.70 84.59	<b>SOTA GNN</b> 77.04 75.93 85.40	<b>P-Value</b> 0.0427 7.79e-22 8.74e-06	<b>P &lt;0.05</b> True True True
Dataset         OGBN-Arxiv         OGBN-Products	Embeddings X-SimTeG X-GIANT X-SimTeG X-GIANT	GraphSAGE           76.84           73.70           84.59           82.84	<b>SOTA GNN</b> 77.04 75.93 85.40 86.12	P-Value 0.0427 7.79e-22 8.74e-06 7.87e-17	P <0.05 True True True True
Dataset OGBN-Arxiv OGBN-Products OGBL Citation2	Embeddings X-SimTeG X-GIANT X-SimTeG X-GIANT X-SimTeG	GraphSAGE           76.84           73.70           84.59           82.84           91.62	<b>SOTA GNN</b> 77.04 75.93 85.40 86.12 91.42	P-Value 0.0427 7.79e-22 8.74e-06 7.87e-17 0.1380	P <0.05 True True True False

Table A8: The results of more LM-involved methods. All results are averaged over 10 runs.

Method	ogbn-arxiv	ogbn-products	ogbl-citation (MRR)
GraphFormer	$72.81 \pm 0.20$	$74.72\pm0.16$	$82.78\pm0.24$
SimTeG + GraphSAGE	$76.84 \pm 0.34$	$84.59\pm0.44$	$85.13\pm0.73$
SimTeG + SOTA GNN	$77.04\pm0.13$	$85.40\pm0.28$	$86.66 \pm 1.21$

Table A9: Statistics of OGBN-Arxiv, OGBN-Products, and OGBL-Citation2-2.7M

Datasets	#Nodes	#Edges	Avg. Degree	#Task	Metric
OGBN-Arxiv (Arxiv)	169,343	1, 166, 243	13.7	node classification	Accuracy
OGBN-Products (Products)	2,449,029	61,859,140	50.5	node classification	Accuracy
OGBL-Citation2-2.7M (Citation2)	2,728,032	27,731,705	10.2	link prediction	MRR / Hits

node's text. For example, "title: multi view metric learning for multi view video summarization; abstract: Traditional methods on video summarization are designed to generate summaries for single-view video records; and thus they cannot fully exploit the redundancy in multi-view video records. In this paper, we present a multi-view metric learning framework for multi-view video summarization that combines the advantages of maximum margin clustering with the disagreement minimization criterion. ..."

OGBN-Products. OGBN-Products is a co-purchase graph, where node denotes a product on Amazon and an edge denotes the co-purchase relationship between two products. The task is to predict the category of each product (node classification). We follow the public split as Hu et al. (2020) and the text processing strategy of GLEM (Zhao et al., 2022). For each node, the corresponding text is its item description. For example, "My Fair Pastry (Good Eats Vol. 9)" "Disc 1: Flour Power (Scones; Shortcakes; Southern Biscuits; Salmon Turnovers; Fruit Tart; Funnel Cake; Sweet or Savory; Pte Choux) Disc 2: Super Sweets 4 (Banana Spitsville; Burned Peach Ice Cream; Chocolate Taffy; Acid Jellies; Peanut Brittle; Chocolate Fudge; Peanut Butter Fudge) ..."

OGBL-Citation2-2.7M. OGBL-Citation2-2.7M is a citation graph, where nodes denote papers and

793

794



Figure A4: Training convergence and validation results of GNNs with X-SimTeG, X-OGB, and X-FIX. The LM backbone is e5-large. The learning rate and batch size are consistent.

edges denote the citations. The task is to predict the missing citation among papers (link prediction). All papers are collected by the official from *Mircrosoft Academic Graph* whereas the text resources are not provided. Though MAG IDs for all papers are provided, we cannot find all corresponding text resources due to the close of MAG project<sup>4</sup>. Hence, we take an intersection of OGBL-Citation2 and OGBN-Papers100M whose text resources are provided by the official, and build a subgraph, namely OGBL-Citation2-2.7M. It contains 93% nodes of OGBL-Citation2 and offers a roughly on-par performance for baselines.

Table A10: The URLs of text resources for ogbn-arxiv, ogbn-products, and OGBL-Citation2.

Dataset	Text Resource URL
OGBN-Arxiv	https://snap.stanford.edu/ogb/data/misc/ogbn_arxiv/titleabs.tsv.gz
OGBN-Products OGBL-Citation2-2.7M	https://drive.google.com/u/0/uc?id=1gsabsx8kkzN9JJZ16J1CA0QA5ASNukhN&export=download https://drive.google.com/u/0/uc?id=19_hkbBUDFZTvQrM0oMbftuXhgz5LbIZY&export=download

#### A2.2 Properties of Language Models

Table A11: Properties of the selected LM backbones. Repositories are hosted by huggingface.

LM	#Params.	#Layers	#Hidden Dim.	Repository
all-MiniLM-L6-v2	23M	6	384	sentence-transformer/all-MiniLM-L6-v2
all-roberta-large-v1	355M	24	1024	sentence-transformer/all-roberta-large-v1
e5-large	335M	24	1024	intfloat/e5-large

#### A2.3 Hyperparameter Search Space

For language models, we design the hyperparameter (HP) search space as in Table A12. Please note that for link prediction, the label smoothing factor is omitted. For HP searching, we utilize optuna (Akiba et al.,

<sup>4</sup>https://www.microsoft.com/en-us/research/project/microsoft-academic-graph/

2019) to search the best HPs for each dataset and each model. For LMs, we take 10 trials. For GNNs, we take 20 trials. The final HP setting for LMs and GNNs are placed as shell scripts in our repository.

	LM			GNN	
hyperparameter	search space	type	hyperparameter	search space	type
learning rate weight decay label smoothing header dropout lora r lora alpha lora dropout	[1e-6, 1e-4] [1e-7, 1e-4] [0.1, 0.7] [0.1, 0.8] [1, 2, 4, 8] [4, 8, 16, 32] [0.1, 0.8]	continual continual continual descrete descrete continual	learning rate weight decay label smoothing dropout num of layers	[1e-4, 1e-2] [1e-7, 1e-4] [0.1, 0.7] [0.1, 0.8] [2, 3, 4, 6, 8]	continual continual continual descrete

Table A12: The search space of LMs and GNNs.