# Making Large Language Models Perform Better in Knowledge Graph Completion

**Anonymous ACL submission** 

#### Abstract

Large language model (LLM) based knowledge 001 graph completion (KGC) aims to predict the missing triples in the KGs with LLMs. However, research about LLM-based KGC fails to sufficiently harness LLMs' inference proficiencies, overlooking critical structural information 007 integral to KGs. In this paper, we explore methods to incorporate structural information into the LLMs, with the overarching goal of facilitating structurally-aware reasoning. We propose a Knowledge Prefix Adapter (KoPA) to fulfill 011 this stated goal. The KoPA uses a structural pre-training phase to comprehend the intricate relations and entities within KGs. Then KoPA communicates such structural understanding to 015 the LLMs through a knowledge prefix adapter 017 which projects the structural embeddings into the textual space and obtains virtual knowledge tokens positioned as a prefix of the input 019 prompt. We conduct comprehensive experiments and provide incisive analysis concerning how the introduction of structural information would be better for LLM's knowledge reasoning ability. Our code and data are available at https://anonymous.4open.science/r/KoPA-0122.

#### 1 Introduction

027

037

041

Knowledge graphs (KGs) (Bollacker et al., 2008) are the quintessential wisdom essence and key infrastructure of modern AI. KGs represent and store real-world knowledge in the triple form: (*head entity, relation, tail entity*). This structured format of knowledge triples offers significant advantages across many AI fields such as recommendation systems (Sun et al., 2020), question answering (Yasunaga et al., 2021), and fault analysis (Chen et al., 2023). However, there is a pertinent drawback of KGs, whether manually curated or automatically extracted. Their scope is restricted to observed knowledge, resulting in an **incomplete** representation riddled with unobserved or missing triples.



Figure 1: A simple case of LLM-based KGC. Useful structural information that describes the surrounding information about the entities can serve as auxiliary prompts and guide the LLM to make correct decisions.

This phenomenon motivates knowledge graph completion (KGC), which aims to predict the missing triples and further enhance the given KG. 042

043

044

045

047

051

053

054

058

059

060

061

062

063

064

065

067

Existing KGC approaches can be divided into two categories: methods based on embeddings (Bordes et al., 2013) and pre-train language models (PLM) (Yao et al., 2019). Recently, as large language models (LLMs) (Zeng et al., 2023; OpenAI, 2023) show outperforming capabilities (Ouyang et al., 2022), this field has recently been revolutionized by LLMs. Some works (Yao et al., 2023) make the first step towards LLM-based KGC, employing existing paradigms like zero-shot reasoning (Brown et al., 2020) and instruction tuning (Ouyang et al., 2022) to accomplish KGC. However, such approaches transform the KGC task into a text-based prediction of individual triples, leading to specific fundamental problems. LLMs lack the depth and precision of factual knowledge which always results in the hallucination (Zhang et al., 2023b) problem of LLMs. Besides, the structural intricacies of KGs such as subgraph structure, relational patterns, and relative entities/relations are often overlooked. This richly structured information, if properly incorporated, can significantly enhance the LLM's understanding and representation of KGs. Figure

068

110

109

111 112 113

114

115

116

117

118 119  1 presents an intuitive view of the importance of structural information for LLM reasoning. However, this is neglected by vanilla IT approaches (Yao et al., 2023) because each input typically only includes a single input triple, leading to potential wastage of the structural information inherent in the KG. Such an approach fails to equip the LLMs with the awareness of the KG structure.

To address these issues, we take a strategic step to LLM-based KGC, aiming to explore how to incorporate the KG structural information into the LLMs and enable structure-aware reasoning. Our initial focus involves transferring the existing LLM paradigms such as in-context learning (ICL) (Dong et al., 2023) and instruction tuning (IT) (Ouyang et al., 2022) to a structure-aware context. We propose a structure-aware ICL method and a structureaware IT method as the base models, focusing on integrating the KG structural information into LLM through text form. Additionally, we propose a Knowledge Prefix Adapter (KoPA) approach to make LLMs a better knowledge reasoner, leveraging structural embedding pre-training to capture the KG structural information. Then KoPA transforms the structural embeddings into textual embedding space by a knowledge prefix adapter and obtains several virtual knowledge tokens. These tokens, acting as prefixes in the input prompt sequence, direct the instruction-tuning process, providing valuable supplementary input triple information. This mapping of structural embeddings to textual form provides auxiliary information to input triples. Besides, we conduct comprehensive analysis and experiments, highlighting the remarkable performance and transferability of KoPA. In summary, our contribution is three-folded:

(1). We are the first extensive investigation of LLM-based KGC, specifically by incorporating KG structural information to enhance the reasoning ability of LLMs. We discuss how to adapt the existing LLM paradigms like ICL and IT to a structure-aware setting for KGC.

(2). We further propose a knowledge prefix adapter (KoPA) that effectively integrates pre-trained KG structural embeddings with LLMs. KoPA fosters a comprehensive interaction between textual embeddings derived from LLMs and structural embeddings sourced from KGs.

(3). We conduct extensive experiments on three public benchmarks and evaluate the KGC performance of all the structure-aware methods proposed by us with adequate baseline comparison with further exploration of the transfer ability and knowledge retention degree.

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

# 2 Related Works

## 2.1 Knowledge Graph Completion

Knowledge graph completion (KGC) (Wang et al., 2017) is an important topic in the KG community, aiming to mine unobserved triples in a given KG. KGC contains several sub-tasks such as triple classification (Bordes et al., 2013), entity prediction (Bordes et al., 2013). The common point among KGC tasks is to establish an effective mechanism to measure the plausibility of the triples. The mainstream KGC methods can be divided into two categories: embedding-based and PLM-based methods. Embedding-based methods (Bordes et al., 2013; Yang et al., 2015; Trouillon et al., 2016; Sun et al., 2019) are designed to embed the entities and relations of KGs into continuous representation spaces. These approaches make full use of structural information from the KGs to model triple plausibility with a well-designed score function and learn the entity/relation embeddings in a self-supervised manner. Moreover, PLM-based methods consider KGC as text-based tasks by fine-tuning pre-trained language models (Devlin et al., 2019). The short textual descriptions are organized as an input sequence and encoded by the PLMs.

## 2.2 LLMs for KG research

Among the research topics of LLM, integrating LLM and KG (Pan et al., 2023) is a popular and important one. On the one hand, hallucination (Zhang et al., 2023b; Yang et al., 2023) is widespread in LLMs which means LLMs are lack factual knowledge and not interpretable. KGs that store structured knowledge can mitigate such a phenomenon (Peng et al., 2023; Feng et al., 2023; Ji et al., 2023) by introducing factual knowledge into LLMs. On the other hand, LLMs can benefit KG-related tasks such as KGC (Zhu et al., 2023b,c), entity alignment (Zhang et al., 2023a), and KGQA (Baek et al., 2023) by its powerful generation capability. KGs for LLMs (KG4LLM) and LLMs for KGs (LLM4KG) are both important research topics. We focus on applying LLMs in the KGC task (LLM4KGC), which has not been carefully studied yet. KGLLaMA (Yao et al., 2023) made the first step by vanilla instruction tuning approach but it lacks in-depth and systematic exploration about how to unleash the power of KGs themselves

219

220

to make structure-aware reasoning in LLMs and
achieve better KGC performance. In this paper, we
will dive into this problem from a more systematic
perspective with the triple classification task.

### **3** Basic Settings for LLM-based KGC

### 3.1 Notations and Preliminaries

173

174

175

176

177

178

179

181

182

183

185

186

188

189

190

191

193

194

195

196

197

198

201

202

206

207

210

211

212

213 214

215

216

218

A KG can be denoted as  $\mathcal{G} = (\mathcal{E}, \mathcal{R}, \mathcal{T}, \mathcal{D})$  where  $\mathcal{E}, \mathcal{R}$  are the entity set, relation set respectively.  $\mathcal{T} = \{(h, r, t) \mid h, t \in \mathcal{E}, r \in \mathcal{R}\}$  is the triple set and  $\mathcal{D}$  is the description set of each entity and relation. We denote  $\mathcal{D}(e), \mathcal{D}(r)$  as the short textual description of each entity  $e \in \mathcal{E}$  and each relation  $r \in \mathcal{R}$ . For example, the text description of the entity '/m/0ctzf1' is  $\mathcal{D}('/m/0ctzf1')$ ='The Transformers'. When applying LLMs to KGC tasks, we denote a LLM as  $\mathcal{M}$  that serves as a text decoder. The input textual sequence S of the model  $\mathcal{M}$  consists of several parts: the instruction prompt  $\mathcal{I}$ , the triple prompt  $\mathcal{X}$ , and the optional auxiliary demonstration prompt  $\mathcal{U}$ . The instruction prompt  $\mathcal{I}$  is the manually prepared instruction to guide the LLM  $\mathcal{M}$  to execute the KGC task. The triple prompt  $\mathcal{X}$  contains the textual information about the triples that need to be processed, which can be denoted as  $\mathcal{X}(h, r, t) = \mathcal{D}(h) \oplus \mathcal{D}(r) \oplus \mathcal{D}(t)$ , where  $(h, r, t) \in \mathcal{T}$  is a triple and  $\oplus$  denotes the textual token concatenation operation. In other words, the short descriptions of h, r, t would be applied as the input information. The auxiliary demonstration prompt  $\mathcal{U}$  is an optional prompt for different settings. In the following, we will follow this set of notations.

Meanwhile, we use triple classification as an entry point to investigate how to utilize LLM to accomplish the KGC task. Triple classification is a basic KGC task aiming to conduct binary classification tasks on the given triples. Whereas in the LLM paradigm, all tasks are converted into the form of text generation. Therefore, we desire the model  $\mathcal{M}$  to answer true or false given the textual sequence input  $S = \mathcal{I} \oplus \mathcal{U} \oplus \mathcal{X}$ .

Triple classification is different from vanilla text classification because the entities and the relations in the prompt have complex semantic information defined by the given KG. Without knowledge of this type of information, the model response is unreliable and unstable. Despite the vast amount of commonsense knowledge that exists in the LLMs (Zhang et al., 2023b), research has shown that large models are numb to fine-grained factual knowledge and will fall into a hallucination. Thus, incorporating the KG information into the prompt to provide more auxiliary information and guide the LLM to make structure-aware reasoning is the key to achieving excellent LLM-based KGC.

#### 3.2 KGC with Existing LLM Paradigms

In this section, we first discuss how to solve the KGC task with existing mainstream LLM paradigms called training-free reasoning approaches and instruction-tuning approaches.

**Training-free reasoning approaches** prompt the LLMs to get direct answers without training. Common training-free methods consist of zero-shot reasoning (ZSR) and in-context learning (ICL). For ZSR, we directly utilize the sequence  $S_{zsr} = \mathcal{I} \oplus \mathcal{X}$ as the input to get the prediction results. For ICL, some demonstration  $\mathcal{U}$  will be added into the input  $S_{icl}$ . To incorporate valuable KG information as the demonstrations, we can sample the relative triples in the local structure of the test triple (h, r, t) to serve as the backbone knowledge. The detailed design of ZSR and ICL approaches are presented in Appendix B.1.1 and B.1.2.

**Instruction tuning approaches** fine-tune the LLMs with instruction template to activate the instruction following ability of LLMs. Vanilla instruction tuning leverages the input  $S_{it}$  to fine-tune LLMs with the next word prediction objective. To incorporate semantic-rich KG information into LLMs, we also propose a structure-aware instruction tuning approach by adding the one-hop neighborhood structure information in the input prompt to inform the LLM with the local structural information. The detailed design of instruction tuning are presented in Appendix B.2.1 and Appendix B.2.2.

Therefore, we provide a detailed discussion of how the existing LLM paradigms can introduce local structural information about KGs to further enhance the model performance. However, though these approaches can work to some extent, they have obvious drawbacks. These **fundamental approaches** to incorporate KG structural information focus on **adding the neighborhood information to the input prompt in the text form**. However, representing the KG structural information in text is not a good choice, which may bring in more invalid or redundant information to the prompt. It's not scalable and effective to increase prompt length indefinitely because a long context will lead to both a decline in model capability and high computa-



Figure 2: An overview of the knowledge prefix adapter (KoPA) by us. KoPA first pre-trains structural-embeddings for the entities and relations in the given KG and then employs instruction tuning to fine-tune the LLM. The structural embeddings of the given input triple will be projected into the textual space of the LLM by the adapter and serve as prefix tokens in the front of the input sequence, which can be "seen" by the following texual tokens due to the unidirectional attention mechanism in the decoder-only LLM.

tional consumption. Besides, we also have difficulty finding the structural information in the KGs that is decisive for triple discrimination. These two problems put us in a dilemma.

#### 4 Methodlogy

271

273

275

276

281

284

287

291

293

294

302

To solve such issues, we propose the Knowledge **P**refix Adapter (KoPA for short) to incorporate the KG structural information into LLM for KGC. Figure 2 presents an intuitive view of KoPA. Firstly we extract the structural information of entities and relations from the KG through structural embedding pre-training, and then we inform this structural information to LLM through a structural prefix adapter into the input sequence S. The LLM Mis further fine-tuned with the structural-enhanced text sequence. We will discuss the details in the next few sections about our design.

#### 4.1 Structural Embedding Pre-training

Instead of adding text about the neighborhood information into the input sequence, KoPA extracts the structural information of the entities and relations by self-supervised structural embedding pretraining. For each entity  $e \in \mathcal{E}$  and each relation  $r \in \mathcal{R}$ , we learn a structural embedding  $e \in$  $\mathbf{R}^{d_e}, \mathbf{r} \in \mathbf{R}^{d_r}$  respectively, where  $d_e, d_r$  are the embedding dimensions. We encode the KG structural information in the embeddings and further adapt them into the textual representation space of LLMs. Referring to the existing embeddingbased KGC paradigm, we define a score function  $\mathcal{F}(h, r, t)$  to measure the plausibility of the triple (h, r, t). We adopt the self-supervised pretraining objective by negative sampling (Bordes et al., 2013):

$$\mathcal{L}_{pre} = \frac{1}{|\mathcal{T}|} \sum_{(h,r,t)\in\mathcal{T}} \left( -\log\sigma(\gamma - \mathcal{F}(h,r,t)) - \sum_{i=1}^{K} p_i \log\sigma(\mathcal{F}(h'_i,r'_i,t'_i) - \gamma) \right)$$
(1)

where  $\gamma$  is the margin,  $\sigma$  is the sigmoid activation function and  $(h'_i, r'_i, t'_i)(i = 1, 2, ..., K)$  are K negative samples (Bordes et al., 2013) of (h, r, t). The weight  $p_i$  is the self-adversarial weights proposed in (Sun et al., 2019).

By minimizing such a pre-training loss, the structural embeddings of each entity and relation are optimized to fit all its relative triples thus the KG structural information such as subgraph structure and relational patterns is captured in the embeddings. Such an approach has been proven effective in many embedding-based KGC methods (Bordes et al., 2013; Sun et al., 2019) to capture classic structural information like relational patterns and distributed entity representations (Hinton et al., 1990) in the earliest days.

#### 4.2 Knowledge Prefix Adapter

After structural embedding pre-training, we could obtain the structural embeddings (h, r, t) of a triple (h, r, t) where the KG structural information is encoded in. However, the structural embeddings are learned in a different representation space against the textual token representation space of the LLM  $\mathcal{M}$ , which means  $\mathcal{M}$  can not directly understand these embeddings. Thus we apply a knowledge prefix adapter  $\mathcal{P}$  to project them into the textual token representation space of  $\mathcal{M}$ . Specifically 304

305

306

307

309

310

311

312

313

314

315

316

317

318

319

320

321

323

324

325

327

328

331

20*1* 

341

344

347

349

367

371

372

speaking, the structural embeddings are converted to several virtual knowledge tokens  $\mathcal{K}$  by  $\mathcal{P}$ :

$$\mathcal{K} = \mathcal{P}(\boldsymbol{h}) \oplus \mathcal{P}(\boldsymbol{r}) \oplus \mathcal{P}(\boldsymbol{t})$$
(2)

In practice, the adapter  $\mathcal{P}$  would be a simple projection layer (Zhu et al., 2023a). Then we put  $\mathcal{K}$ in the front of the original input sequence S serving as a prefix of the instruction and triple prompt  $\mathcal{S}_{kpa} = \mathcal{K} \oplus \mathcal{I}_{it} \oplus \mathcal{X}$ . This way, all the following text tokens can be seen with the prefix  $\mathcal{K}$  due to the unidirectional attention in decoder-only LLMs. By doing this, the textual tokens can pay unidirectional attention to the structural embeddings of the input triple. Such a structure-aware prompt will be employed during fine-tuning and inference. During training, we froze the pre-trained structural embeddings. The adapter is optimized to learn the mapping from structural knowledge toward textual representation and will have the generalization to new triples in the inference stage, which will benefit the textual description and provide the triple information from another perspective to make enhanced predictions.

### 4.3 Complexity Analysis

After proposing KoPA, we make a comparison among LLM-based KGC methods to demonstrate the advantages of KoPA, which is shown in Table 3. Compared with the basic paradigms (ZSR/ICL/IT), KoPA incorporates the KG structural embeddings into LLM to combine the textual and structural information. Meanwhile, KoPA makes the length of the prompt more refined as the length of virtual tokens generated by the structural prefix adapter is fixed to 3 for head/relation/tail respectively. In contrast, the prompt length of structure-aware IT (enhanced IT in the table) is linearly related to the number of neighborhood triples k. In contrast to methods that incorporate structural information based on textual descriptions, KoPA achieves this goal by fixed-length virtual knowledge tokens generated by the adapter.

## 5 Experiments

## 5.1 Datasets

In our experiments, we use three public KG benchmarks UMLS (Yao et al., 2019), CoDeX-S (Safavi and Koutra, 2020), and FB15K-237N (Lv et al., 2022) to evaluate the proposed LLM-based KGC methods. The detailed split information of the datasets is shown in Table 4 of the Appendix.

## 5.2 Experimental Settings

## 5.2.1 Baseline Methods

In our experiments, we provide a comprehensive comparison with three broad classes of baselines on triple classification, which is an important subtask of KGC. The KGC baselines can be divided into three parts: embedding-based methods (Bordes et al., 2013; Yang et al., 2015; Trouillon et al., 2016; Sun et al., 2019), PLM-based methods (Yao et al., 2019; Lv et al., 2022), and LLM-based methods (Yao et al., 2023). Besides, we further divide the LLM-based methods into two categories: trainingfree methods and fine-tuning methods. Trainingfree methods consist of ZSR and ICL, while finetuning methods consist of vanilla IT and structureaware IT (enhanced IT). 381

382

383

384

387

388

389

391

392

393

394

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

## 5.2.2 Implementation and Detail Settings

We reproduce the baseline results and implement the KoPA proposed by us. We employ Alpaca-7B (Taori et al., 2023) as the LLM backbone. Alpaca is a famous extended version of LLaMA (Touvron et al., 2023) model fine-tuned on instructionfollowing data. We reproduce the triple classification results of KGLLaMA (Yao et al., 2023) over two backbones (LLaMA and Alpaca) to avoid the effect of backbone choice on the results. We name the two baseline models KGLLaMA and KGAlpaca respectively. For all the fine-tuning methods (instruction tuning, structure-aware instruction tuning, and KoPA), we fine-tune Alpaca using LoRA (Hu et al., 2022) with rank 64. The number of epochs is searched in  $\{3, 4, 5\}$  and the learning rate is tuned in  $\{1e^{-4}, 3e^{-4}, 5e^{-4}\}$ . We use the AdamW optimizer (Loshchilov and Hutter, 2019) with a fixed batch size of 12. We conducted all the experiments with Nvidia A800 GPUs. The embedding pre-training process is efficient which only takes several minutes. We make a detailed discussion about the time cost of experiments in Appendix C.4.

## 5.2.3 Evaluation Protocol

We evaluate the methods with triple classification task (Bordes et al., 2013), which is essentially binary classification and all the test datasets are labelbalanced. Therefore, we use accuracy, precision, recall, and F1-score as the evaluation metrics.

## 5.3 Main Results

The main experiment results of triple classification are shown in Table 1. Since precision and

Table 1: The main experiment results of triple classification. We report the accuracy (ACC), precision (P), recall (R), and F1-score (F1) results for each method on the three datasets. "-" means the result are missing because the specificity of PKGC makes it difficult to reproduce. The best **Acc / F1** results in baselines are marked with <u>underline</u>, and we highlight our results with bold when we achieve new SOTA.

	Model	Model		UMLS		CoDeX-S			FB15K-237N				
	mouth	Acc	P	R	F1	Acc	P	R	F1	Acc	Р	R	F1
	TransE (Bordes et al., 2013)	84.49	86.53	81.69	84.04	72.07	71.91	72.42	72.17	69.71	70.80	67.11	68.91
Embadding based	DistMult (Yang et al., 2015)	86.38	87.06	86.53	86.79	66.79	69.67	59.46	64.16	58.66	58.98	56.84	57.90
Enibedding-based	ComplEx (Trouillon et al., 2016)	90.77	89.92	91.83	90.87	67.64	67.84	67.06	67.45	65.70	66.46	63.38	64.88
	RotatE (Sun et al., 2019)	<u>92.05</u>	90.17	94.41	<u>92.23</u>	75.68	75.66	75.71	75.69	68.46	69.24	66.41	67.80
DI M based	KG-BERT (Yao et al., 2019)	77.30	70.96	92.43	80.28	77.30	70.96	92.43	80.28	56.02	53.47	97.62	67.84
PLW-based	PKGC (Lv et al., 2022)	-	-	-	-	-	-	-	-	<u>79.60</u>	-	-	79.50
	Zero-shot(Alpaca)	52.64	51.55	87.69	64.91	50.62	50.31	99.83	66.91	56.06	53.32	97.37	68.91
	Zero-shot(GPT-3.5)	67.58	88.04	40.71	55.67	54.68	69.13	16.94	27.21	60.15	86.62	24.01	37.59
LLM-based	ICL(1-shot)	50.37	50.25	75.34	60.29	49.86	49.86	50.59	50.17	54.54	53.67	66.35	59.34
Training-free	ICL(2-shot)	53.78	52.47	80.18	63.43	52.95	51.54	98.85	67.75	57.81	56.22	70.56	62.58
	ICL(4-shot)	53.18	52.26	73.22	60.99	51.14	50.58	99.83	67.14	59.29	57.49	71.37	63.68
	ICL(8-shot)	55.52	55.85	52.65	54.21	50.62	50.31	99.83	66.91	59.23	57.23	73.02	64.17
	KG-LLaMA (Yao et al., 2023)	85.77	87.84	83.05	85.38	79.43	78.67	80.74	79.69	74.81	67.37	96.23	79.25
LLM-based	KG-Alpaca (Yao et al., 2023)	86.01	94.91	76.10	84.46	80.25	79.38	81.73	80.54	69.91	62.71	98.28	76.56
Fine-tuning	Vanilla IT	86.91	95.18	77.76	85.59	81.18	77.01	88.89	82.52	73.50	65.87	97.53	78.63
-	Structure-aware IT	89.93	93.27	86.08	89.54	<u>81.27</u>	77.14	88.40	<u>82.58</u>	76.42	69.56	93.95	<u>79.94</u>
KoPA		92.58	90.85	94.70	92.70	82.74	77.91	91.41	84.11	77.65	70.81	94.09	80.81

recall alone do not give a good response to the model's performance on the classification task, we focus on accuracy and F1-score. However, to provide a comprehensive analysis of different models, we also report the precision and recall results in the table. Overall, we can find that KoPA achieves outperforming accuracy and F1 results compared with the existing 16 baseline models on all three datasets. Taking CoDeX-S as an example, KoPA achieves 1.81% improvement in accuracy and 1.85% improvement on F1. As we use the pre-trained RotatE embeddings in KoPA, we can observe that KoPA significantly outperforms the original embedding-based RotatE method, especially on larger and more challenging datasets like CoDeX-S and FB15K-237N.

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

Meanwhile, compared with all LLM-based approaches, we can see that the LLMs cannot understand the KG structural information well without fine-tuning. The zero-shot LLMs perform very poorly in the triple classification task even though GPT-3.5-turbo (175B parameters) has excellent capability. Though the demonstrations provided by ICL can incorporate the KG information, the performance gain is limited. Besides, the prediction results of training-free methods are biased and easy to slip into the extremes of all-right or all-wrong, as the recall of them is either very high or very low but the F1 scores are relatively low all the time.

However, fine-tuning LLMs can introduce the KG information into LLMs as the overall per-

formance makes obvious improvements. Meanwhile, though structure-aware IT enhances the input prompt with neighborhood information of triples, its performance is also limited compared with KoPA. This suggests that the structural embeddings consist of more semantic-rich information compared with text-based auxiliary prompts, which can also be understood by the LLM through the prefix adapter. Combining the analysis in Section 4.3 and the experimental results, KoPA achieves better results on top of shorter prompts.

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

### 5.4 Transferability Exploration

The results in the main experiments have shown the effectiveness of KoPA. To further validate the generality and the transferability of KoPA, we conduct a new transferability experiment. In this experiment, we will demonstrate that the knowledge prefix adapter will learn to transfer from structural embeddings to textual token representations and provide semantic-rich auxiliary information to enhance the decoding process of LLM inference.

We demonstrate this point by testing the influence of KoPA for entities that do not appear in the training phase, which is also called inductive setting in other KGC works (Chen et al., 2022b). We split the KG dataset into an inductive setting with a defined inductive rate (IR), which refers to the ratio of unseen entities during training. For example, if IR=10%, we will randomly select 10% entities as the inductive entity set. Any triple in the training



Figure 3: The results of the transferbility experiment. We report the results on CoDeX-S dataset under different inductive rate (IR). Besides, we split the test data into seen (S) and unseen (U) parts based on whether the entity appeared during training. Also we total the results of all (A) the test data together. Accuracy (Acc) and F1-score (F1) are reported in the radar charts.

491

492

493

494

495

496

497

498

499

502

505

506

507

510

513

514

515

517

set whose head or tail is in the inductive set will be removed during training. Besides, the triples in the test set will be divided into two parts: the seen (S) part and the unseen (U) part. If the head or tail in a triple is in the inductive entity set, it will be regarded as unseen. We fine-tune the LLM with only remaining seen triples and test on both seen and unseen triples. In this setting, a set of entities will not participate in the training process and the LLM does not see their textual descriptions, which will make the test process more challenging. We report the accuracy and F1 score for seen (S), unseen (U), and all (A) test triples, which is shown in Figure 3 for three fine-tuning methods: KoPA, vanilla IT, and structure-aware IT (enhanced IT in the figure).

From the radio charts, we can observe that KoPA outperforms the other methods for unseen triples and has less performance degradation when IR increases. The performance of structure-aware IT (enhanced IT) with neighborhood triples in the textual form is more unstable. These phenomena suggest that the knowledge prefix adapter can learn a good mapping from the structural embeddings to the textual representation, which is transferable even if the entities are unseen during training. The structural embeddings captured from KG play a more significant role in informing the LLM with useful structural information.

Table 2: Ablation study results on CoDeX-S. We first replace the pre-trained structural embedding with other components and change the insert position of virtual knowledge tokens to demonstrate the effectiveness of knowledge prefix adapter.

M	Acc	F1	
KoPA(Pref	82.74	84.11	
Embedding	w/o SE	81.18	82.52
	w/ TransE	82.46	83.42
	w/ DistMult	80.71	81.27
	w/ ComplEx	81.21	82.12
	w/ Random	81.53	82.36
Position	Infix	81.21	82.69
	Suffix	77.29	77.75

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

553

554

555

### 5.5 Ablation Study

To verify the effectiveness of the KoPA design, we conduct a two-part ablation study. The first part is designed to verify the effectiveness of structural embedding and the second part is designed to verify the effectiveness of prefix adapter. As shown in Table 2, we can find that removing the structural embeddings or replacing them with random initialized embeddings both lead to performance decline. Also, we find that the model is compatible with different types of structural embeddings. However, the performance gain depends on whether the embedding was originally powerful in the triple classification task or not. Refer to Tables 1, TransE (Bordes et al., 2013) and RotatE (Sun et al., 2019) are better embedding-based KGC models compared with DistMult (Yang et al., 2015) and ComplEx (Trouillon et al., 2016). This demonstrates that semantic-rich structural information is the key to performance improvement and KoPA takes full advantage of it.

Meanwhile, putting the virtual knowledge tokens generated by the adapter in the middle (infix) or in the last (suffix) of the input sequence will also decrease the performance. We believe the reason is that putting tokens in the front of the sequence will make all the text pay attention to them as LLMs are usually decoder-only architectures with unidirectional self-attention. Then the LLM can make a better decision with the structural embeddings that fully interact with the text. Combining these two parts of the ablation study, we believe that our design of KoPA is effective and reasonable.

### 5.6 Case Study

To make a more intuitive view of KoPA, we conduct a case study in this section from both macro and micro perspectives. From a macro perspective, we count the prediction overlap of several models and



Figure 4: The Venn diagram of the correct predictions from various KGC models. Each intersecting part in the diagram represents the same predictions from different models on certain data.

plot a Venn diagram shown in Figure 4.

556

561

562

564

566

567

571

573

574

576

579

581

583

588

589

591

594

From the diagram we can find that KoPA has a significant portion of the proper predictions that do not intersect with several other models, which means that KoPA makes the right prediction on some test data that many other models predict incorrectly. This suggests that the structural information incorporated in KoPA has a significant role in making correct predictions. For a micro example, a test triple (John Landis, film director film, Coming to America) is predicted as wrong by the RotatE model and vanilla instruction tuning LLM. With retrieved neighborhood triples (Coming to America, locations, New York City), (John Landis, nationality, USA ), (Coming to America, genre, romantic comedy), (Comedy, common netflix titles, Coming to America), the structure-aware fine-tuned LLM still makes a wrong prediction because the neighborhood information is of little use in the judgment of the current prediction though they are the correct factual. The structural embeddings applied in KoPA contain more information than structural information in the form of text and are easier for us to extract by a structural pre-training process. Thus, KoPA outperforms other models in the triple classification task.

#### 5.7 Common Ability Retention

To delve into the preservation of generic capabilities in LLMs, we conducted another experiment to assess the overall proficiency of LLMs both before and after fine-tuning. We apply the MMLU (Hendrycks et al., 2021) benchmark for this problem. MMLU is the most popular benchmark to evaluate the general abilities of LLMs in different domains such as Humanities, Social Sciences, STEM, and others. The overall evaluation results on different datasets are shown in Figure 5:

From the results, it can be noticed that after KoPA training, there were discernible alterations



Figure 5: The common ability experiments on MMLU.

595

596

597

598

599

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

in the generalized abilities of LLMs. In most instances, there was a decrease, but notably, STEM proficiency exhibited improvement on the UMLS dataset. We attribute this phenomenon to the UMLS being a medical KG, encompassing substantial knowledge in medicine, biology, and chemistry, and training on this dataset allows the model to acquire more STEM knowledge. Consequently, when facing natural language inputs differing from the training task, the model adeptly leverages the acquired knowledge from KGC task fine-tuning to get enhanced results. We have listed several subjects in MMLUs that showed improvement after training with UMLS. These subjects are highly relevant and close to the knowledge domain encapsulated in the UMLS. the LLMs trained with the KGC task also achieved significant improvements across different input prompts, marking a compelling observation.

#### 6 Conclusion

In this paper, we systematically explore how to incorporate structural information into LLMs to make structure-aware reasoning for KGC tasks. We extend the original LLM paradigms and propose structure-aware ICL and IT methods to incorporate the structural information by text. We further propose KoPA, a knowledge prefix adapter to incorporate the pre-trained structural embeddings into the LLMs. We conduct triple classification experiments to make comprehensive comparisons among the structure-aware methods and demonstrate the outperforming results achieved by KoPA. In the future, we plan to dive deep into LLM-based KGC and think about a more unified framework to accomplish all the KGC tasks with LLMs. Besides, we will explore adapting KGs into LLMbased downstream applications to make the LLMs knowledgeable about KGs.

732

733

734

735

736

## 632 Limitations

- In this paper, we focus on the topic of LLM-based
  KGC, aiming to integrate the KGC capability into
  LLM with structure-aware reasoning ability. There
  are some limitations in our work.
- 637 Model Design. The design of the knowledge prefix
  638 adapter is relatively simple. A more sophisticated
  639 and subtle design will be our future plan.
- Task Generalization. we have not generalized
  the model method to all kinds of KGC tasks such
  as link prediction. The main task of the current
  research is triple classification.
- 644 Model Interpretability. More in-depth experi645 ments need to be further conducted to explore the
  646 interpretability of our model..

### Ethical Considerations

647

651

658

667

670

671

674

675

676

677 678

679

The data and pre-trained LLMs we use in our paper are all existing open-sourced resources for academic research and there are no ethical issues in our paper. We promise that there is no anti-ethic behavior in our research.

#### References

- Jinheon Baek, Alham Fikri Aji, and Amir Saffari. 2023. Knowledge-augmented language model prompting for zero-shot knowledge graph question answering. *CoRR*, abs/2306.04136.
- Kurt D. Bollacker, Colin Evans, Praveen K. Paritosh, Tim Sturge, and Jamie Taylor. 2008. Freebase: a collaboratively created graph database for structuring human knowledge. In *SIGMOD Conference*, pages 1247–1250. ACM.
- Antoine Bordes, Nicolas Usunier, Alberto García-Durán, Jason Weston, and Oksana Yakhnenko. 2013. Translating embeddings for modeling multirelational data. In *NIPS*, pages 2787–2795.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *NeurIPS*.
- Chen Chen, Yufei Wang, Bing Li, and Kwok-Yan Lam. 2022a. Knowledge is flat: A seq2seq generative framework for various knowledge graph completion.

In *COLING*, pages 4005–4017. International Committee on Computational Linguistics.

- Mingyang Chen, Wen Zhang, Yushan Zhu, Hongting Zhou, Zonggang Yuan, Changliang Xu, and Huajun Chen. 2022b. Meta-knowledge transfer for inductive knowledge graph embedding. In *SIGIR*, pages 927– 937. ACM.
- Zhuo Chen, Wen Zhang, Yufeng Huang, Mingyang Chen, Yuxia Geng, Hongtao Yu, Zhen Bi, Yichi Zhang, Zhen Yao, Wenting Song, Xinliang Wu, Yi Yang, Mingyi Chen, Zhaoyang Lian, Yingying Li, Lei Cheng, and Huajun Chen. 2023. Tele-knowledge pre-training for fault analysis. In *ICDE*, pages 3453– 3466. IEEE.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT (1)*, pages 4171–4186. Association for Computational Linguistics.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, Lei Li, and Zhifang Sui. 2023. A survey for in-context learning. *CoRR*, abs/2301.00234.
- Chao Feng, Xinyu Zhang, and Zichu Fei. 2023. Knowledge solver: Teaching llms to search for domain knowledge from knowledge graphs. *CoRR*, abs/2309.03118.
- Jiayan Guo, Lun Du, and Hengyu Liu. 2023. Gpt4graph: Can large language models understand graph structured data ? an empirical evaluation and benchmarking. *CoRR*, abs/2305.15066.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. In *ICLR*. OpenReview.net.
- Geoffrey E. Hinton, James L. McClelland, and David E. Rumelhart. 1990. Distributed representations. In *The Philosophy of Artificial Intelligence*, Oxford readings in philosophy, pages 248–280. Oxford University Press.
- 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 *ICLR*. OpenReview.net.
- Ziwei Ji, Zihan Liu, Nayeon Lee, Tiezheng Yu, Bryan Wilie, Min Zeng, and Pascale Fung. 2023. RHO: reducing hallucination in open-domain dialogues with knowledge grounding. In ACL (Findings), pages 4504–4522. Association for Computational Linguistics.
- Bosung Kim, Taesuk Hong, Youngjoong Ko, and Jungyun Seo. 2020. Multi-task learning for knowledge graph completion with pre-trained language models. In *COLING*, pages 1737–1743. International Committee on Computational Linguistics.

737

738

- 770 773 774 775 778 779
- 781
- 783 784

- Youwei Liang, Ruiyi Zhang, Li Zhang, and Pengtao Xie. 2023. Drugchat: Towards enabling chatgptlike capabilities on drug molecule graphs. CoRR, abs/2309.03907.
- Chang Liu and Bo Wu. 2023. Evaluating large language models on graphs: Performance insights and comparative analysis. CoRR, abs/2308.11224.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. CoRR. abs/2304.08485.
- Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In ICLR (Poster). Open-Review.net.
- Xin Lv, Yankai Lin, Yixin Cao, Lei Hou, Juanzi Li, Zhiyuan Liu, Peng Li, and Jie Zhou. 2022. Do pretrained models benefit knowledge graph completion? A reliable evaluation and a reasonable approach. In ACL (Findings), pages 3570–3581. Association for Computational Linguistics.
- Chenyang Lyu, Minghao Wu, Longyue Wang, Xinting Huang, Bingshuai Liu, Zefeng Du, Shuming Shi, and Zhaopeng Tu. 2023. Macaw-llm: Multi-modal language modeling with image, audio, video, and text integration. CoRR, abs/2306.09093.
- OpenAI. 2023. GPT-4 technical report. CoRR, abs/2303.08774.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In NeurIPS.
- Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Jiapu Wang, and Xindong Wu. 2023. Unifying large language models and knowledge graphs: A roadmap. CoRR, abs/2306.08302.
- Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng, Yujia Xie, Yu Hu, Qiuyuan Huang, Lars Liden, Zhou Yu, Weizhu Chen, and Jianfeng Gao. 2023. Check your facts and try again: Improving large language models with external knowledge and automated feedback. CoRR, abs/2302.12813.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yangi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. J. Mach. Learn. Res., 21:140:1-140:67.
- Tara Safavi and Danai Koutra. 2020. Codex: A comprehensive knowledge graph completion benchmark. In EMNLP (1), pages 8328-8350. Association for Computational Linguistics.

Apoorv Saxena, Adrian Kochsiek, and Rainer Gemulla. 2022. Sequence-to-sequence knowledge graph completion and question answering. In ACL (1), pages 2814–2828. Association for Computational Linguistics.

790

791

792

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

- Ming Shi and Jing Zhao. 2022. A knowledge graph link prediction model with combined 1d and 2d convolutional embeddings. In DSAA, pages 1-2. IEEE.
- Rui Sun, Xuezhi Cao, Yan Zhao, Junchen Wan, Kun Zhou, Fuzheng Zhang, Zhongyuan Wang, and Kai Zheng. 2020. Multi-modal knowledge graphs for recommender systems. In CIKM, pages 1405–1414. ACM.
- Zhiqing Sun, Zhi-Hong Deng, Jian-Yun Nie, and Jian Tang. 2019. Rotate: Knowledge graph embedding by relational rotation in complex space. In ICLR (Poster). OpenReview.net.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https:// github.com/tatsu-lab/stanford\_alpaca.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models. CoRR, abs/2302.13971.
- Théo Trouillon, Johannes Welbl, Sebastian Riedel, Éric Gaussier, and Guillaume Bouchard. 2016. Complex embeddings for simple link prediction. In ICML, volume 48 of JMLR Workshop and Conference Proceedings, pages 2071–2080. JMLR.org.
- Denny Vrandecic and Markus Krötzsch. 2014. Wikidata: a free collaborative knowledgebase. Commun. ACM, 57(10):78-85.
- Bo Wang, Tao Shen, Guodong Long, Tianyi Zhou, Ying Wang, and Yi Chang. 2021. Structure-augmented text representation learning for efficient knowledge graph completion. In WWW, pages 1737-1748. ACM / IW3C2.
- Heng Wang, Shangbin Feng, Tianxing He, Zhaoxuan Tan, Xiaochuang Han, and Yulia Tsvetkov. 2023a. Can language models solve graph problems in natural language? CoRR, abs/2305.10037.
- Quan Wang, Zhendong Mao, Bin Wang, and Li Guo. 2017. Knowledge graph embedding: A survey of approaches and applications. IEEE Trans. Knowl. Data Eng., 29(12):2724–2743.
- Yufei Wang, Wanjun Zhong, Liangyou Li, Fei Mi, Xingshan Zeng, Wenyong Huang, Lifeng Shang, Xin Jiang, and Qun Liu. 2023b. Aligning large language models with human: A survey. CoRR, abs/2307.12966.

942

943

944

945

946

947

948

949

950

- 847 848
- 850 851
- 85 85
- 8
- 8
- 8
- 8
- 8
- 8
- 867
- 869 870
- 8

871

888 889 890

887

89

89 89

897 898 898

- Bishan Yang, Wen-tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. 2015. Embedding entities and relations for learning and inference in knowledge bases. In *ICLR (Poster)*.
- Linyao Yang, Hongyang Chen, Zhao Li, Xiao Ding, and Xindong Wu. 2023. Chatgpt is not enough: Enhancing large language models with knowledge graphs for fact-aware language modeling. *CoRR*, abs/2306.11489.
- Liang Yao, Chengsheng Mao, and Yuan Luo. 2019. KG-BERT: BERT for knowledge graph completion. *CoRR*, abs/1909.03193.
- Liang Yao, Jiazhen Peng, Chengsheng Mao, and Yuan Luo. 2023. Exploring large language models for knowledge graph completion. *CoRR*, abs/2308.13916.
- Michihiro Yasunaga, Hongyu Ren, Antoine Bosselut, Percy Liang, and Jure Leskovec. 2021. QA-GNN: reasoning with language models and knowledge graphs for question answering. In *NAACL-HLT*, pages 535–546. Association for Computational Linguistics.
- Hongbin Ye, Ningyu Zhang, Hui Chen, and Huajun Chen. 2022. Generative knowledge graph construction: A review. In *EMNLP*, pages 1–17. Association for Computational Linguistics.
- Ruosong Ye, Caiqi Zhang, Runhui Wang, Shuyuan Xu, and Yongfeng Zhang. 2023. Natural language is all a graph needs. *CoRR*, abs/2308.07134.
- Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. 2023. A survey on multimodal large language models. *CoRR*, abs/2306.13549.
- Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, Weng Lam Tam, Zixuan Ma, Yufei Xue, Jidong Zhai, Wenguang Chen, Zhiyuan Liu, Peng Zhang, Yuxiao Dong, and Jie Tang. 2023. GLM-130B: an open bilingual pre-trained model. In *ICLR*. OpenReview.net.
- Rui Zhang, Yixin Su, Bayu Distiawan Trisedya, Xiaoyan Zhao, Min Yang, Hong Cheng, and Jianzhong Qi. 2023a. Autoalign: Fully automatic and effective knowledge graph alignment enabled by large language models. *CoRR*, abs/2307.11772.
- Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. 2023b. Siren's song in the AI ocean: A survey on hallucination in large language models. *CoRR*, abs/2309.01219.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao

Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023. A survey of large language models. *CoRR*, abs/2303.18223.

- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023a. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *CoRR*, abs/2304.10592.
- Yuqi Zhu, Xiaohan Wang, Jing Chen, Shuofei Qiao, Yixin Ou, Yunzhi Yao, Shumin Deng, Huajun Chen, and Ningyu Zhang. 2023b. Llms for knowledge graph construction and reasoning: Recent capabilities and future opportunities. *CoRR*, abs/2305.13168.
- Yuqi Zhu, Xiaohan Wang, Jing Chen, Shuofei Qiao, Yixin Ou, Yunzhi Yao, Shumin Deng, Huajun Chen, and Ningyu Zhang. 2023c. Llms for knowledge graph construction and reasoning: Recent capabilities and future opportunities. *CoRR*, abs/2305.13168.
- Yuqi Zhu, Xiaohan Wang, Jing Chen, Shuofei Qiao, Yixin Ou, Yunzhi Yao, Shumin Deng, Huajun Chen, and Ningyu Zhang. 2023d. Llms for knowledge graph construction and reasoning: Recent capabilities and future opportunities. *CoRR*, abs/2305.13168.

# Appendix

# **A Related Works**

In this section, we add some related works about KG and LLM that for reasons of space could not be mentioned in the main paper.

# A.1 Knowledge Graph Completion

Depending on the design of the scoring function, embedding-based methods can be divided into three sub-categories: (1) translation-based methods like TransE (Bordes et al., 2013) and RotatE (Sun et al., 2019), (2) tensor decomposition methods like DistMult (Yang et al., 2015) and ComplEx (Trouillon et al., 2016), (3) neural networkbased methods like ConvE (Shi and Zhao, 2022). Embedding-based KGC methods learn the structural embeddings for triple discrimination but neglect the textual information in the KG.

KG-BERT (Yao et al., 2019) is the first PLMbased method that models KGC as a binary text classification task. Subsequent works like MTL-KGC (Kim et al., 2020) and StAR (Wang et al., 2021) have further improved KG-BERT by introducing more training tasks such as relation classification and triple ranking and more complex triple encoding strategy. PKGC (Lv et al., 2022) utilizes manual prompt templates to capture the triple semantic. Other methods like KGT5 (Saxena et al., 2022) and KG-S2S (Chen et al., 2022a) make a

1004

1006

1007

1008

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1027

1028

1029

1031

1032

1033

1034

1035

1036

1037

1038

1040

1041

1042

step on the generative KGC (Ye et al., 2022) in a sequence-to-sequence paradigm with encoderdecoder PLMs like T5 (Raffel et al., 2020). PLMbased methods leverage the power of PLM but make the training process into text-based learning, which is difficult to capture complex structure information in the KGs.

### A.2 LLMs for KG research

951

952

953

955

956

957

960

961

962

963

964

966

968

969

970

972

973

974

975

976

977 978

979

982

983

984

987

996

997

1000

In recent years, large language models (LLMs) (OpenAI, 2023; Zeng et al., 2023; Touvron et al., 2023) have made rapid progress and demonstrated powerful capabilities in a considerable number of text-related tasks (Zhao et al., 2023). LLMs are usually pre-trained in an auto-regressive manner with next word prediction task (Brown et al., 2020) and demonstrate strong capability on text comprehension and generation. Some significant techniques such as instruction tuning (IT) (Ouyang et al., 2022) and human preference alignment (Wang et al., 2023b) are further applied to guide the model to follow human instructions and generate responses that are consistent with human values and preferences.

### A.3 Incorporate Non-textual Modality Information into LLMs

As LLMs demonstrate generalizable capabilities on text generation, many other works attempt to incorporate non-textual modality such as images (Liu et al., 2023; Zhu et al., 2023a), audio (Lyu et al., 2023), and video (Lyu et al., 2023), which are also called multi-modal LLMs (Yin et al., 2023). These methods tend to encode non-textual information through the modality encoders and then process it as virtual text tokens. The non-textual tokens are aligned with the word tokens by instruction tuning on multi-modal datasets.

The multi-modal LLM mentioned above usually excludes graph, which is another important data modality. There are also some works talking about how to incorporate graph data into LLMs. Drug-Chat (Liang et al., 2023) proposes to encode the drug molecule graphs with graph encoders and finetune the LLM to predict drug interactions. Other works (Ye et al., 2023; Liu and Wu, 2023; Wang et al., 2023a; Guo et al., 2023) explore how to solve graph learning tasks like node classification and graph classification by convert the graph structure information into LLMs.

Our research is relative to this topic as KGs also have complex graph structures on top of the text descriptions. In this paper, we will explore how to incorporate complex structural information in the KGs into the LLMs to achieve better reasoning capabilities on knowledge graph completion.

## B Base Paradigms for LLMs to Solve KGC tasks

### **B.1** Training-free Reasoning Approaches

Training-free reasoning is an efficient approach to employing a LLM to solve downstream tasks without extra training. We need to prepare a suitable prompt template to acquire the results generated by the model  $\mathcal{M}$ . Mainstream training-free approaches consist of zero-shot reasoning and incontext learning (Dong et al., 2023). Existing methods like (Zhu et al., 2023d) have tried zero-shot reasoning to evaluate the link prediction ability of LLM. Besides, there are no ICL-based KGC methods yet. We will discuss each of them below more systematically and incorporate structural information into LLMs.

#### **B.1.1 Zero-shot Reasoning Approach**

Zero-shot reasoning (ZSR) is a direct approach for LLMs to do the reasoning task without auxiliary information  $\mathcal{U}$ . Thus, the input sequence of ZSR can be denoted as  $S_{zsr} = \mathcal{I}_{zsr} \oplus \mathcal{X}$ . The decoding process of the LLM  $\mathcal{M}$  can be formulated as:

$$\mathcal{A}_{zsr} = \arg \max_{\mathcal{A}} P_{\mathcal{M}}(\mathcal{A}|\mathcal{S}_{zsr})$$
  
=  $\arg \max_{\mathcal{A}} P_{\mathcal{M}}(\mathcal{A}|\mathcal{I}_{zsr}, \mathcal{X})$  (3)

where  $\mathcal{A}$  is the generated answer of the model  $\mathcal{M}$ and  $\mathcal{I}_{zsr}$  is the instruction template for ZSR. In the setting of ZSR, no KG information is added to the input sequence  $\mathcal{S}_{zsr}$ .

The determinative information in the ZSR prompt is only the textual descriptions of the test triple. ZSR is unable to incorporate KG information due to its setting limitations, otherwise, it cannot be called zero-shot.

### B.1.2 In-context Learning Approach with Structure-aware Demonstration

As another training-free paradigm, in-context learning (ICL) (Dong et al., 2023) allows the model  $\mathcal{M}$ to add auxiliary demonstration  $\mathcal{U}$  to the input Sand accomplish the task in the form of analogical reasoning, which can be denoted as:

$$\mathcal{A}_{icl} = \arg \max_{\mathcal{A}} P_{\mathcal{M}}(\mathcal{A}|\mathcal{S}_{icl})$$
  
=  $\arg \max_{\mathcal{A}} P_{\mathcal{M}}(\mathcal{A}|\mathcal{I}_{icl},\mathcal{U},\mathcal{X})$  (4) 1043

1086

1087

1088

1090

As for the triple classification task, the demonstration  $\mathcal{U}$  should be some triples and their labels in the form of  $\{(\mathcal{X}_i, y_i), 1 \leq i \leq k\}$ , where  $\mathcal{X}_i$  is the demonstration triple and  $\dagger_i$  is the label. We denote the ICL with k demonstrations as k-shot ICL.

The demonstration triples can be randomly sampled from the existing training KG. However, to further incorporate the relative KG information of the test triple (h, r, t), we propose to sample triples that are in the local structure of h and t, which means one of the entities in each sampled triple should be h or t. Besides, as existing KG only consists of positive triples, we employ negative sampling (Lv et al., 2022) to sample negative triples for demonstration. The number of positive and negative triples are the same for balanced predictions. In the demonstration prompt, the positive triples are labeled as true and the negative triples are labeled as false.

By doing this, we incorporate the local structural information into the demonstration prompt  $\mathcal{U}$  with both positive and negative samples. Such a structure-aware demonstration could better enhance the analogical reasoning process of the model  $\mathcal{M}$ .

#### B.2 Instruction Tuning Approaches

Instruction tuning (IT) aims to fine-tune the LLM to follow human instructions and accomplish the mentioned tasks in the instruction prompt. In this section, we will talk about how to incorporate the KG information into IT approaches.

#### **B.2.1** Vanilla Instruction Tuning

In the setting of vanilla IT, the instruction prompt  $\mathcal{I}_{it}$  will describe the details of completing the triple classification task and the triple prompt  $\mathcal{X}$  consists of the input triple. No other auxiliary demonstrations are included in the input template. To train the model  $\mathcal{M}$ , the input sequence is organized as  $\mathcal{S}_{it} = \mathcal{I}_{it} \oplus \mathcal{X} \oplus \mathcal{A}_{it}$ . where  $\mathcal{A}_{it}$  is the predicted answer of the training data. The model  $\mathcal{M}$  is finetuned with the next word prediction task (Zhao et al., 2023) which is a universal approach to training LLMs. The training objective can be formulated as:

$$\mathcal{L}_{it} = -\frac{1}{|\mathcal{S}_{it}|} \sum_{i=1}^{|\mathcal{S}_{it}|} \log P_{\mathcal{M}}(s_i|s_{< i}) \tag{5}$$

where  $s_i(i = 1, 2, ..., |S_{it}|)$  represents the textual tokens of the input sequence  $S_{it}$ . In the inference

Table 3: Comparasion among LLM-based KGC methods in three ways. As for the prompt length anaysis,  $L_I$ ,  $L_T$  denote the length of the instruction prompt and triple prompt.  $L_D$  denotes the length of a demonstration and k is the demonstration number. ZSR/ICL/IT refer to zero-shot reasoning, in-context learning, and instruction tuning respectively.

Method	Requires Fine-tuning	Extra KG Info	Prompt Length
ZSR	×	×	$L_I + L_T$
ICL	×	$\checkmark$	$L_I + L_T + kL_D$
Vanilla IT	$\checkmark$	×	$L_I + L_T$
Enhanced IT	$\checkmark$	$\checkmark$	$L_I + L_T + kL_D$
KoPA	$\checkmark$	$\checkmark$	$L_I + L_T + 3$

stage, the model  $\mathcal{M}$  is employed to predict the answer  $\mathcal{A}_{it}$  of the test data like Equation 3. Besides, negative sampling (Lv et al., 2022) is also applied as training KG only consists of positve triples.

1091

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

Vanilla IT only fine-tunes the LLM to learn the knowledge in the single triple to discriminate. Such an approach makes it difficult to fully utilize the rich semantics present in a KG and the model performance is limited.

#### **B.2.2** Structure-aware Instruction Tuning

As mentioned before, the structural information of KG plays a significant role in the KGC tasks (Wang et al., 2017). To incorporate such KG information during the fine-tuning stage, we achieve this goal by adding the neighborhood descriptions of the input triple. Specifically, we can sample the neighborhoods of the head h and tail t and put the textual descriptions of neighborhood triples in the demonstration prompt  $U_{it}$ . In this way, the input training sequence is enhanced as  $S_{it} = I_{it} \oplus U_{it} \oplus X \oplus A_{it}$ .

We name such an approach as structure-aware instruction tuning as the local structural information of the entities is added into the input sequence in the form of neighborhood triples.

### B.3 Comparasions Among Baselines and KoPA

We make a comparasion table of the base paradigms and KoPA, presented in Table 4.3.

#### **C** Experiments

### C.1 Dataset Details

UMLS (Yao et al., 2019) is a classic medical knowl-<br/>edge graph including general knowledge about<br/>medicine and health care. CoDeX-S (Safavi and<br/>Koutra, 2020) is an encyclopedic KG extracted1121<br/>1122

Table 4: Statistical information of datasets. The positve (+) and negative (-) samples are 1:1 in the valid and test set.

Dataset	$ \mathcal{E} $	$ \mathcal{R} $	#Train	#Valid(+/-)	#Test(+/-)
UMLS	135	46	5216	652/652	661/661
CoDeX-S	2034	42	32888	1827/1827	1828/1828

from Wikidata (Vrandecic and Krötzsch, 2014). FB15K-237N proposed in (Lv et al., 2022) is modified from FB15K-237. Besides, CoDeX-S and FB15K-237N mine hard negative triples for a more challenging evaluation and avoid false negative samples in the validation/test dataset during dataset construction. We constructed negative samples for UMLS in the same method.

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

Our dataset selection is based on three key perspectives:

(1). **Dataset prevalence**. The three KG datasets in our experiments are all mainstream public datasets, which are widely acknowledged and extensively employed in KG-related research.

(2). **Dataset size**. We intentionally varied the dataset sizes, ranging from small (UMLS) to medium (CoDeX-S) to large (FB15K-237N), to explore the generalizability of our approach across datasets of different scales.

(3). **Dataset diversity**. Our choice also takes into account the diversity of the dataset. UMLS is a domain-specific KG dataset with substantial medical knowledge, while CoDeX-S and FB15K-237N are then two encyclopaedic KG datasets, offering a balance between general world knowledge and specialized domains.

From the experimental results, our method shows some improvement in all three datasets, which indicates the generality of our method.

### C.2 Baseline Details

The specific models used for these baselines are listed below:

(1). Embedding-based KGC methods. We select four traditional embedding-based KGC methods for comparisons, namely TransE (), DistMult (Yang et al., 2015), ComplEx (Trouillon et al., 2016), and RotatE (Sun et al., 2019). These methods predict the triple plausibility by the learned structural embeddings and the score functions defined in the model.

(2). **PLM-based KGC methods**. We select KG-BERT (Yao et al., 2019) and PKGC (Lv et al.,

2022) as PLM-based KGC baselines, which are<br/>classic methods focusing on the triple classification1167task. These methods treat triple classification as a<br/>binary text classification task.1169

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1205

1206

(3). **LLM-based KGC methods**. LLM-based KGC research is still at an early stage. There are only KGLLaMA (Yao et al., 2023) to be the LLM-based KGC baseline. In addition to KGLLaMA, the methods proposed in Section 3 by us including ZSR, ICL, IT, and structure-aware IT (enhanced IT) will also serve as baselines.

#### C.3 Implementation Details for Baselines

For embedding-based KGC methods, we reproduce the results with OpenKE we set the embedding dimension  $d_e = d_r = 512$  and sample K = 32negative samples during training. The margin  $\gamma$  is tuned among {0, 4, 6, 8, 12}. After training KGC models, we search for the best classification score threshold on the validation set for test data following the traditional setting (Bordes et al., 2013).

For PLM-based methods, the backbone model for PLM-based KGC methods is BERT (Devlin et al., 2019). We fine-tune the KG-BERT according to the official code implementation. Since PKGC requires a lot of manual work to annotate each relation with a prompt, we only report the results of FB15K-237N shown in the original paper.

For zero-shot reasoning, in addition to measuring with the same backbone Alpaca, we also test the performance of the *GPT-3.5-turbor* which has 175B parameters. For the in-context learning method, we sample k-shot (k=1,2,4,8) structureaware demonstrations. Besides, we sample 4 neighborhood triples for each triple to conduct structureaware instruction tuning. For KoPA, we employ RotatE (Sun et al., 2019) and the score function of structural embedding pre-training and the embedding dimension is set to 512 and the adapter is a  $512 \times 4096$  linear projection layer.

### C.4 Training Efficiency

We would like to emphasize the notable distinction 1207 between structural embedding pre-training and the 1208 pre-training methods applied in language models 1209 (such as BERT, GPT, and LLaMA), as mentioned 1210 in our paper. Language model pre-training is con-1211 ducted on the massive corpus in a self-supervised 1212 manner with transformer models which is computa-1213 tionally expensive. However, structural embedding 1214 pre-training is performed on a certain KG and the 1215

	w/o Training	w/ Training
Clinical	44.9	47.9
College Medicine	30.1	31.2
High School Biology	42.9	46.8
High School Chemistry	30.0	32.0
Medical Genetics	44.0	48.0

Table 5: The specific domains in MMLU which LLM achieves higher scores after training on UMLS.

embedding parameters and complexity is relatively modest compared with language models.

For each individual experiment, we use one A800 GPU with 80G of memory for both training and inference. Our linux server has ubuntu installed and three A800 GPUs are available, this is our detailed equipment situation. To provide more precise insight, our pre-training process takes only about ten minutes on a single A800 GPU, whereas the LLM fine-tuning process takes serveral hours for various datasets (1 hour for UMLS, 3-4 hours for CoDeX-S, 8-9 hours for FB15K-237N). Pretraining a language model like BERT and GPT requires dozens of GPUs and days of training. The discernible contrast in efficiency between these two pre-training methodologies underscores the lightweight nature of structural embedding pretraining.

## C.5 Common Ability Study

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

We list some of the specific domains that LLMshave improved after training on the UMLS asshown in Table 5.