

Beyond Prefixes: Graph-as-Memory Cross-Attention for Knowledge Graph Completion with Large Language Models

Anonymous ACL submission

Abstract

Fusing Knowledge Graphs with Large Language Models (LLMs) is crucial for knowledge-intensive tasks like knowledge graph completion. Existing LLM-based approaches typically inject graph information via prefix concatenation, resulting in shallow interactions that fail to support fine-grained evidence retrieval during generation. Beyond prefixes, we propose Graph-as-Memory Tuning (GMT), a new paradigm that represents local graph structure as explicit graph memory and injects it into LLMs via deep, token-wise cross-attention. Specifically, GMT first employs a Semantic Graph Module to encode context-aware semantics from local neighborhoods guided by knowledge-enhanced relations, and compresses them into a fixed number of graph memory tokens. A Graph-as-Memory Cross-Attention Fusion Module then integrates these tokens into multiple Transformer layers, allowing LLM hidden state to dynamically retrieve relevant graph evidence. To enable efficient adaptation, GMT applies LoRA only to the memory cross-attention while keeping the base LLM frozen. Extensive experiments show that GMT significantly outperforms prefix-tuning and other strong baselines, providing more potent signals for robust reasoning.

1 Introduction

Knowledge Graphs (KGs) organize complex facts into structured triples (h, r, t) , empowering applications ranging from recommendation systems (Cui et al., 2025; Chen et al., 2025) to question answering (Omar et al., 2023; Lu et al., 2025). However, their inherent incompleteness necessitates Knowledge Graph Completion (KGC) to infer missing links from existing facts (Pan et al., 2024; Li et al., 2025a).

Traditionally, KGC has relied on *embedding-based models*, including geometric approaches (Bordes et al., 2013; Sun et al., 2019) and

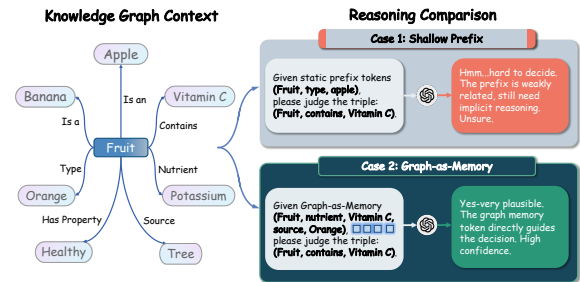


Figure 1: The guiding effect of structured information on LLMs. The semantics of the relation “Treats” change dynamically based on the graph context.

graph neural networks (Dettmers et al., 2018; Zhu et al., 2021). While proficient at encoding static structural patterns, these methods often overlook the rich textual semantics inherent in entities and relations. To address this, recent research has pivoted towards *LLM-based paradigms* (Yao et al., 2019; Li et al., 2024; Zhang et al., 2024a), which leverage the generative and semantic capabilities of pre-trained models. However, existing integration strategies, predominantly based on prefix-tuning, typically employ a shallow fusion approach that simply concatenates structural embeddings with textual inputs (Zhang et al., 2024a). This shallow interaction fails to deeply align structural signals with textual representations, imposing a heavy implicit reasoning burden on the LLM and often leading to hallucinations or context-insensitive predictions.

Consequently, a critical challenge remains: *how to effectively fuse explicit KG structure with implicit LLM semantics at a deep, feature-interactive level?* As illustrated in Figure 1, relational semantics are dynamic and context-dependent. The relation “Treats” implies distinct mechanisms, such as “Symptomatic Relief” or “Pathogen Targeted Treatment”, depending entirely on the local graph neighborhood (e.g., *Aspirin* vs. *Oseltamivir*). Capturing these nuanced shifts requires more than static concatenation; it demands a mechanism that can

072 dynamically modulate the LLM’s perception based
073 on structural context.

074 To bridge this gap, we propose **Graph-as-**
075 **Memory Tuning (GMT)**, a memory-centric frame-
076 work that reframes local graph structure as an ex-
077 plicit graph memory and integrates it into LLMs
078 through deep, token-wise cross-attention. GMT
079 comprises two core components: 1) A **Semantic**
080 **Graph Module** that uses a relation centric mes-
081 sage passing mechanism to extract a dense and
082 context-aware semantics from the local neighbor-
083 hood, guided by knowledge enhanced relation de-
084 scriptions, and compresses them into a fixed num-
085 ber of graph memory tokens. 2) A **Graph-as-**
086 **Memory Cross-Attention Fusion Module** that
087 injects these graph memory tokens into multiple
088 Transformer layers, enabling each prompt token
089 hidden representation to dynamically retrieve rele-
090 vant evidence from the graph memory during gen-
091 eration. To maintain parameter efficiency, GMT
092 keeps the base LLM frozen and applies LoRA only
093 to the memory cross-attention pathway. Our main
094 contributions are summarized as follows:

- 095 • We propose **GMT**, a deep fusion paradigm that
096 replaces shallow concatenation with memory-
097 based, token-wise retrieval via cross-attention,
098 bridging graph structure and LLM semantics.
099 (Section 3)
- 100 • We introduce a **Semantic Graph Module** that
101 leverages knowledge-enhanced relation seman-
102 tics to guide neighborhood aggregation and con-
103 struct compact graph memory tokens. (Section
104 3.2)
- 105 • We design a **Graph-as-Memory Cross-**
106 **Attention Fusion Module** that performs
107 multi-layer memory injection and token-wise
108 retrieval, and enable parameter-efficient training
109 by applying LoRA-based adaptation to align
110 graph memory with a frozen LLM. (Section 3.3)
- 111 • Empirical results confirm that GMT achieves
112 state-of-the-art performance on multiple KGC
113 benchmarks, validating the efficacy of deep in-
114 jection. (Section 4)

115 2 Related Work

116 2.1 Knowledge Graph Completion

117 Traditional KGC methods predominantly rely on
118 embedding-based paradigms to capture structural
119 patterns. Geometric models, such as TransE (Bor-
120 des et al., 2013) and RotatE (Sun et al., 2019),
121 map entities to continuous spaces using transla-

122 tional or rotational scoring functions. Extensions
123 like HAKE (Zhang et al., 2020) and BoxE (Ab-
124 boud et al., 2020) further incorporate hierarchical
125 and logical constraints. Parallely, GNN-based ap-
126 proaches like ConvE (Dettmers et al., 2018) and
127 NBFNet (Zhu et al., 2021) capture local topology
128 through message passing. Despite their structural
129 efficacy, these methods assign static representations
130 to entities, limiting their ability to generalize to un-
131 seen data or leverage the rich semantics inherent in
132 KGs (Chang et al., 2024).

133 2.2 LLMs for KGC

134 The adoption of Large Language Models has signif-
135 icantly advanced semantic reasoning in KGC. Early
136 LLM-based discriminative approaches, such as KG-
137 BERT (Yao et al., 2019) and SimKGC (Wang et al.,
138 2022), formulate KGC as a sequence classifica-
139 tion task but do not explicitly model graph topol-
140 ogy. In the generative setting, instruction-tuned
141 frameworks including KICGPT (Wei et al., 2023)
142 and KG-LLaMA (Zhu and De Meo, 2025) directly
143 predict tail entities, while later methods such as
144 MKGL (Guo et al., 2024b) and KG-FIT (Jiang
145 et al., 2024) incorporate multi-view learning and
146 few-shot inductive reasoning capabilities. More re-
147 cent frameworks seek to integrate explicit structural
148 information with the semantic reasoning capabili-
149 ties of LLMs. Representative approaches include
150 KoPA (Zhang et al., 2024a), which project struc-
151 tural embeddings into the textual space as prefix
152 tokens; SSQR (Lin et al., 2025), which optimizes
153 this via quantized representations; and GLTW (Luo
154 et al., 2025), which employs a joint graph trans-
155 former architecture. Nevertheless, these methods
156 largely rely on shallow integration strategies such
157 as prefix concatenation or textualization, failing to
158 establish a deep, feature-level interaction in which
159 graph structure dynamically modulates the internal
160 representations of the LLM.

161 3 Methodology

162 3.1 Preliminaries

163 A knowledge graph \mathcal{G} is defined as a set of triples
164 $\mathcal{T} = \{(h, r, t) \mid h, t \in \mathcal{E}, r \in \mathcal{R}\}$, where \mathcal{E} and
165 \mathcal{R} denote the entity and relation sets, respectively.
166 Knowledge Graph Completion (KGC) aims to in-
167 fer a missing element in an incomplete triple, e.g.,
168 predicting t for a query $(h, r, ?)$. In LLM-based
169 KGC, a model \mathcal{M} is prompted with a textualized
170 query and trained to generate (or select) the missing

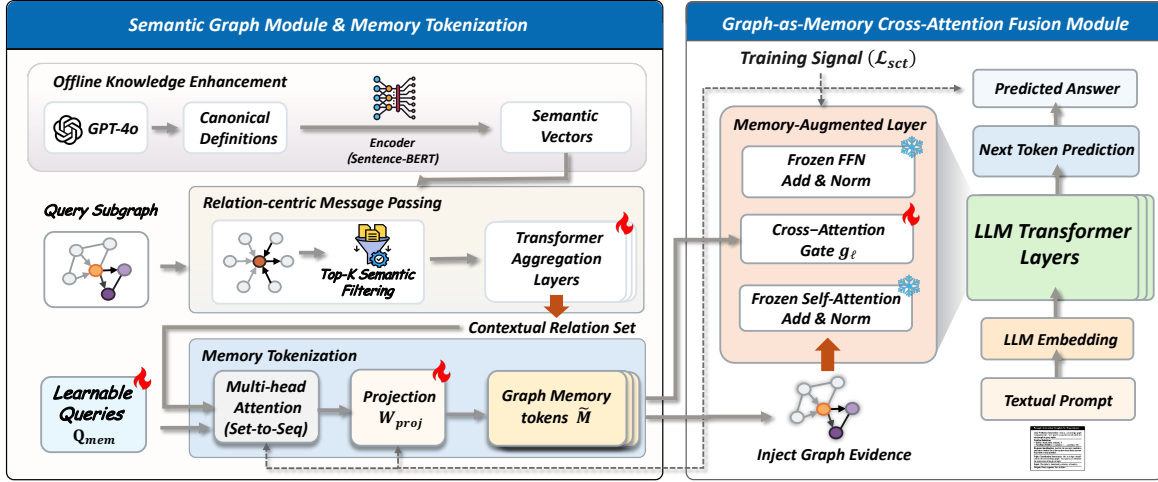


Figure 2: The GMT framework.

entity.

As shown in Figure 2, our GMT conditions the LLM on a graph memory built from the local semantic subgraph around the query entities, and injects such memory into multiple Transformer layers through cross-attention, enabling deep and token-wise retrieval of graph evidence during generation.

3.2 Semantic Graph Module

The first stage of GMT is to transform the local neighborhood structure into a dense set of contextual semantic representations that can serve as external memory for the LLM. Unlike methods that directly rely on pre-trained entity embeddings (Zhang et al., 2024a), our Semantic Graph Module (SGM) performs *relation-centric* message passing to extract semantically filtered relational evidence around the entities.

Relation-centric Message Passing. Inspired by (Wang et al., 2021; Li et al., 2025b), we treat relations as the primary carriers of KG semantics¹. For a given query triple (h, r, t) , we extract the local neighborhoods around h and t (masking the predicted entity for link prediction). For each central edge e_c incident to h or t , we aggregate information from its neighboring edges $e_n \in \mathcal{N}(e_c)$. To mitigate noise from indiscriminate aggregation, we perform Top- K neighbor filtering using explicit semantic relevance. Specifically, we first conduct an offline **Knowledge Enhancement** step for each re-

¹For instance, the rule $(A, \text{is_father_of}, B) \wedge (C, \text{is_wife_of}, A) \rightarrow (C, \text{is_mother_of}, B)$ is resolved by relational interplay. Moreover, an entity’s identity is often implicitly encoded by its relational context, making relation-centric aggregation effective.

lation type using a strong LLM (GPT-4o²), producing canonical definitions (Prompt in Appendix A), and encode them with an embedding model (e.g., Sentence-BERT³ (Reimers and Gurevych, 2019)) to obtain a semantic vector for each relation. During message passing, for each central edge e_c and a neighbor edge e_n , we compute relevance by cosine similarity between their pre-computed semantic vectors s_c and s_n :

$$\text{Score}(e_c, e_n) = \frac{s_c \cdot s_n}{\|s_c\|_2 \|s_n\|_2}. \quad (1)$$

We then select $\mathcal{N}_K(e_c)$ as the Top- K neighbors with highest scores and aggregate their states (e.g., mean pooling) to obtain $\bar{s}^{\mathcal{N}_K(e_c)}$.

We refine each central edge representation via a Transformer-style update, where s_c^l attends to the aggregated neighborhood context:

$$s_c^{l+1} = \text{TransformerLayer}(s_c^l, \bar{s}^{\mathcal{N}_K(e_c)}), \quad (2)$$

and iterate for L layers to obtain contextual edge representations around h and t .

Graph Memory Tokenization. A single pooled vector inevitably bottlenecks diverse evidence. Instead, we expose a **set of memory tokens** to the LLM. Let $\mathcal{S} = \{s_{h,i}^L\}_i \cup \{s_{t,j}^L\}_j \in \mathbb{R}^{N \times d_c}$ denote the set of contextual relation states collected from both sides (with N variable per query). We compress \mathcal{S} into a fixed-length graph memory with m tokens using a learnable set-to-sequence tokenizer. Concretely, we introduce m learnable memory queries $\mathbf{Q}_{mem} \in \mathbb{R}^{m \times d_c}$ and compute:

$$\mathbf{M} = \text{Attn}(\mathbf{Q}_{mem}, \mathcal{S}, \mathcal{S}) \in \mathbb{R}^{m \times d_c}, \quad (3)$$

²<https://github.com/openai/openai-python>.

³<https://github.com/UKPLab/sentence-transformers>

where $\text{Attn}(\cdot)$ is standard multi-head attention. \mathbf{M} serves as a compact yet expressive **Graph-as-Memory** representation for the subsequent fusion module. Finally, we project \mathbf{M} to the LLM hidden size:

$$\tilde{\mathbf{M}} = \mathbf{M}\mathbf{W}_{proj}, \quad \mathbf{W}_{proj} \in \mathbb{R}^{d_c \times d_{llm}}. \quad (4)$$

3.3 Graph-as-Memory Cross-Attention Fusion

The core challenge is to inject the structured graph evidence into the LLM in a *deep* and *token-wise* manner. We propose a Graph-as-Memory Cross-Attention Fusion Module that allows each token to retrieve relevant KG evidence from $\tilde{\mathbf{M}}$ at multiple layers.

Memory-augmented Transformer Layers. Given the input prompt embeddings $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, the LLM produces hidden states $\mathbf{H}^0 = \mathbf{X}$ and updates them through L_{llm} Transformer layers. For a subset of layers $\mathcal{L}_{mem} \subseteq \{1, \dots, L_{llm}\}$, we insert a cross-attention sub-layer after self-attention:

$$\hat{\mathbf{H}}^\ell = \mathbf{H}^\ell + \text{SelfAttn}(\text{LN}(\mathbf{H}^\ell)), \quad (5)$$

$$\mathbf{H}^{\ell+1} = \hat{\mathbf{H}}^\ell + g_\ell \cdot \text{CrossAttn}(\hat{\mathbf{H}}^\ell, \tilde{\mathbf{M}}), \quad (6)$$

followed by the standard FFN block. Here $\text{CrossAttn}(\cdot)$ uses queries from token states and keys/values from the graph memory, enabling each token to selectively retrieve graph evidence. We use a learnable gate g_ℓ for training stability, allowing the model to gradually incorporate memory.

Low-Rank Adaptation of Cross-Attention. To adapt the memory-reading pathway efficiently, we apply LoRA to the projection matrices of cross-attention (and only these parameters, unless otherwise stated). For each projection matrix $\mathbf{W} \in \{\mathbf{W}_q, \mathbf{W}_k, \mathbf{W}_v, \mathbf{W}_o\}$ in the cross-attention module, LoRA parameterizes:

$$\mathbf{W} = \mathbf{W}_0 + \Delta\mathbf{W}, \quad \Delta\mathbf{W} = \mathbf{B}\mathbf{A}, \quad (7)$$

where \mathbf{W}_0 is frozen, $\mathbf{A} \in \mathbb{R}^{r \times d}$, $\mathbf{B} \in \mathbb{R}^{d \times r}$, and $r \ll d$. Only \mathbf{A}, \mathbf{B} are trained. This design keeps the base LLM frozen while opening a dedicated, parameter-efficient channel to align graph memory with the LLM latent space.

At inference, we construct graph memory $\tilde{\mathbf{M}}$ from the query neighborhood and generate the missing entity autoregressively:

$$\mathcal{A} = \arg \max_{\mathcal{A}} P_{\mathcal{M}}(\mathcal{A} | \mathcal{I}_{GMT}, \mathbf{X}, \tilde{\mathbf{M}}), \quad (8)$$

where \mathcal{I}_{GMT} is the instruction template and \mathcal{A} is the predicted answer.

3.4 Training Strategy

We adopt a two-stage training paradigm. Stage 1 pre-trains the graph module to capture structural/semantic regularities, providing a strong initialization. Stage 2 aligns the graph memory and the LLM through memory-augmented fine-tuning with LoRA on cross-attention.

3.4.1 Stage 1: Self-Supervised Pre-training

We first pre-train SGM on a self-supervised link prediction objective. Given (h, r, t) , the graph module produces contextualized representations for the head and tail, denoted as \mathbf{e}_h and \mathbf{e}_t (obtained by mean pooling the final relation states associated with each entity). We define a plausibility score:

$$F_{\mathcal{G}}(h, r, t) = \langle \mathbf{e}_h, \mathbf{e}_r, \mathbf{e}_t \rangle, \quad (9)$$

where \mathbf{e}_r is the knowledge-enhanced relation semantic vector (from GPT-4o and Sentence-BERT). With negative sampling, for each positive triple $(h, r, t) \in \mathcal{T}$, we construct corrupted triples $\mathcal{T}' = \{(h'_i, r, t'_i)\}_{i=1}^k$. The loss is:

$$\mathcal{L}_{\text{Graph}} = - \sum_{(h,r,t) \in \mathcal{T}} \left[\log \sigma(F_{\mathcal{G}}(h, r, t)) + \sum_{i=1}^k \log \sigma(-F_{\mathcal{G}}(h'_i, r, t'_i)) \right]. \quad (10)$$

This stage equips SGM with robust relational semantics before interacting with the LLM.

3.4.2 Stage 2: Memory-Augmented Alignment with the LLM

Starting from the pre-trained SGM, we fine-tune the full GMT pipeline on the KGC objective while keeping the base LLM frozen. For each training query, we build the graph memory $\tilde{\mathbf{M}}$ using Eq. (3) and feed the textual prompt into the LLM equipped with memory cross-attention layers (Eq. (6)). We optimize the auto-regressive next-token prediction loss:

$$\mathcal{L}_{GMT} = - \sum_{i=1}^{|\mathcal{S}_{GMT}|} \log P(s_i | s_{<i}), \quad (11)$$

where $\mathcal{S}_{GMT} = \mathcal{I}_{GMT} \oplus \mathbf{X} \oplus \mathcal{A}$. Gradients from \mathcal{L}_{GMT} update: (i) the graph memory tokenizer parameters and projection \mathbf{W}_{proj} , and (ii) the LoRA weights of cross-attention in the selected layers. This training aligns the graph memory space with

Algorithm 1 Memory-Augmented Fine-tuning Stage of GMT

- 1: **Input:** Knowledge graph \mathcal{G} ; training triple (h, r, t) ; pre-trained SGM_θ ; memory tokenizer parameters ω (including \mathbf{Q}_{mem}); memory projection \mathbf{W}_{proj} ; frozen base LLM; cross-attn LoRA parameters $\psi_{mem-lora}$.
 - 2: **Output:** Updated $(\omega, \mathbf{W}_{proj}, \psi_{mem-lora})$.
 - 3: Prepare textual input:
 - 4: $S_{in} \leftarrow \text{Textualize}(h, r)$;
 - 5: $\mathbf{X} \leftarrow \text{Embedding}_{LLM}(S_{in})$;
 - 6: Build graph memory (mask the predicted entity for link prediction):
 - 7: $\mathcal{S} \leftarrow SGM_\theta(h, t, \mathcal{G})$; $\{\mathcal{S}$: contextual relation set
 - 8: $\tilde{\mathbf{M}} \leftarrow \text{Attn}(\mathbf{Q}_{mem}, \mathcal{S}, \mathcal{S})$;
 - 9: $\tilde{\mathbf{M}} \leftarrow \mathbf{M}\mathbf{W}_{proj}$;
 - 10: Memory-augmented forward with LoRA cross-attention:
 - 11: $P(\cdot | \mathbf{X}, \tilde{\mathbf{M}}) \leftarrow LLM(\mathbf{X}; \tilde{\mathbf{M}}, \psi_{mem-lora})$;
 - 12: Compute loss and update parameters:
 - 13: $\mathcal{L}_{GMT} \leftarrow -\sum_{j=1}^{|t|} \log P(t_j | \mathbf{X}, \tilde{\mathbf{M}}, t_{<j})$;
 - 14: Update $(\omega, \mathbf{W}_{proj}, \psi_{mem-lora})$ by backprop;
 - 15: **Return** updated parameters.
-

Table 1: Statistics of the benchmark datasets. For triple classification datasets (UMLS, CoDeX-S, FB15k-237N), the validation/test splits are denoted as positive/negative samples.

Dataset	Task	#Entities	#Relations	#Train	#Valid	#Test
WN18RR	Link	40,943	11	86,835	3,034	3,134
FB15k-237	Link	14,541	237	272,115	17,535	20,466
UMLS	Triple	135	46	5216	652/652	661/661
CoDeX-S	Triple	2034	42	32888	1827/1827	1828/1828
FB15k-237N	Triple	13,104	93	87,282	7041/7041	8226/8226

the LLM’s internal representations through a dedicated memory-reading pathway, enabling deep fusion without modifying the base model weights. The detailed workflow is presented in Algorithm 1.

4 Experiments

To evaluate the effectiveness of GMT, we designed a series of experiments to address the following research questions:

- **RQ1:** How significantly can GMT enhance LLMs’ performance in KGC tasks?
- **RQ2:** What are the distinct contributions of GMT’s key components to its overall performance?
- **RQ3:** How Does Context from Semantic Graphs Enhance the Reasoning of LLMs?

4.1 Experiment Settings

Datasets. We evaluate our GMT framework on various widely recognized KGC benchmark datasets:

WN18RR (Dettmers et al., 2018) and **FB15k-237** (Toutanova and Chen, 2015): These datasets

are used for the **link prediction** task. Both have inverse relations removed, making them standards for evaluating complex reasoning abilities.

UMLS, CoDeX-S and **FB15k-237N** (Lv et al., 2022): These are employed for the **triple classification** task. It provides high-quality negative samples for each triple, making them ideal for assessing a model’s ability to distinguish factual correctness. Detailed statistics are provided in Table 1.

Baselines. Our selection of baselines covers a wide spectrum of representative methods, from traditional KGE models to the latest LLM-based approaches. Many results are directly cited from the SSQR (Lin et al., 2025), Kopa (Zhang et al., 2024a) and GLTW (Luo et al., 2025) to maintain consistency.

For Link Prediction. We compare against (1) *Emb-based Methods*: TransE (Bordes et al., 2013), CompGCN (Vashishth et al., 2020), AdaProp (Zhang et al., 2024b), MA-GNN (Xu et al., 2023), TCRA (Guo et al., 2024a), ARR (Chen et al., 2024), and DiffusionE (Cao et al., 2024); (2) *LLM-based Methods*: KICGPT (Wei et al., 2023), CSProm-KG-CD (Li et al., 2023), KG-FIT (Jiang et al., 2024), MKGL (Guo et al., 2024b), SSQR-LLaMA2 (Lin et al., 2025).

For Triple Classification. We compare against (1) *Emb-based Methods*: TransE (Bordes et al., 2013), DistMult (Yang et al., 2015), ComplexE (Trouillon et al., 2016), and RotatE (Sun et al., 2019); (2) *LLM-based Methods*: KG-LLaMA (Zhu and De Meo, 2025), KG-Alpaca (Zhu and De Meo, 2025), KOPA (Zhang et al., 2024a), SSQR-LLaMA2 (Lin et al., 2025) and SAT (Liu et al., 2025).

Implementation Details. We build GMT on Alpaca-7B with the base LLM frozen. We pre-train the SGM for 20 epochs with self-supervised link prediction using negative sampling ($k=64$), and then fine-tune GMT with the base LLM frozen. We use $K=8$ for Top- K neighbor filtering and $m=32$ memory tokens, injecting memory via cross-attention in the top layers $\mathcal{L}_{mem}=\{25, \dots, 32\}$ with a learnable gate initialized to zero. In Stage 2, we train only the memory tokenizer, the memory projection, and LoRA on cross-attention projections $\{W_q, W_k, W_v, W_o\}$ (rank $r=64$, $\alpha=128$, dropout 0.01). We tune epochs in $\{3,4,5\}$ and learning rate in $\{1e-4, 3e-4, 5e-4\}$, using AdamW with batch size 12, bf16, and gradient clipping 100.0. Experiments are run on Nvidia A800-80GB

Table 2: Link Prediction on FB15k-237 and WN18RR datasets. Best results are in **bold**, second best are underlined.

Type	Model	WN18RR				FB15k-237			
		MRR	Hits@1	Hits@3	Hits@10	MRR	Hits@1	Hits@3	Hits@10
Emb-based	TransE	0.223	0.014	0.401	0.529	0.330	0.231	0.369	0.528
	CompGCN	0.479	0.443	0.494	0.546	0.355	0.264	0.390	0.535
	AdaProp	0.562	0.499	–	0.671	0.417	0.331	–	0.585
	MA-GNN	0.565	0.507	0.592	0.679	0.379	0.282	0.415	0.569
	TCRA	0.496	0.457	0.511	0.574	0.367	0.275	0.403	0.554
	ARR	0.521	–	0.607	–	0.398	–	0.436	–
LLM-based	DiffusionE	0.557	0.504	–	0.658	0.376	0.294	–	0.539
	KICGPT	0.549	0.474	0.585	0.641	0.412	0.327	0.448	0.554
	CSProm-KG-CD	0.559	0.508	0.578	0.660	–	–	–	–
	KG-FIT	0.553	0.488	0.595	0.695	0.362	0.275	0.485	0.572
	MKGL	0.552	0.500	0.577	0.656	0.415	0.325	0.454	0.591
	SSQR-LLaMA2	0.591	0.548	0.618	0.673	0.449	0.374	0.491	0.597
Ours	GLTW _{7b}	0.593	0.556	0.649	0.690	0.469	0.351	0.481	0.614
Ours	GMT	0.621	0.569	0.667	0.703	0.488	0.394	0.505	0.629

Table 3: Triple classification on UMLS, CoDeX-S, and FB15K-237N datasets. Best results are in **bold**, second best are underlined.

Type	Model	UMLS				CoDeX-S				FB15K-237N			
		Acc	P	R	F1	Acc	P	R	F1	Acc	P	R	F1
Emb-based	TransE	84.49	86.53	81.69	84.04	72.07	71.91	72.42	72.17	69.71	70.80	67.11	68.91
	DistMult	86.38	87.06	86.53	86.79	66.79	69.67	59.46	64.16	58.66	58.98	56.84	57.90
	ComplEx	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	92.05	90.17	94.41	92.23	75.68	75.66	75.71	75.69	68.46	69.24	66.41	67.80
LLM-based	Alpaca _{0-shot}	52.64	51.55	87.69	64.91	50.62	50.31	99.83	66.91	56.06	53.32	97.37	68.91
	GPT-3.5 _{0-shot}	67.58	88.04	40.71	55.67	54.68	69.13	16.94	27.21	60.15	86.62	24.01	37.59
	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	85.77	87.84	83.05	85.38	79.43	78.67	80.74	79.69	74.81	67.37	96.23	79.25
	KG-Alpaca	86.01	94.91	76.10	84.46	80.25	79.38	81.73	80.54	69.91	62.71	98.28	76.56
	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
	SAT	92.24	91.05	93.99	93.17	85.55	83.38	89.31	86.54	82.71	82.30	85.24	83.28
Ours	GMT	94.55	91.86	95.74	93.76	89.01	83.27	92.43	87.61	84.10	83.14	91.27	87.02

GPUs. All closed-source models used for knowledge enhancement are accessed via their official APIs, using default inference hyperparameters (e.g., temperature) with no additional tuning.

4.2 Overall Performance Comparison: RQ1

Link Prediction Classic link prediction is a ranking task that evaluates the position of the gold entity among candidates. We fine-tune GMT with an instruction template \mathcal{I}_{GMT} (Appendix A) and follow SSQR-LLaMA2’s candidate setting by using a pre-trained AdaProp to retrieve 20 candidates per query. As shown in Table 2, GMT achieves state-of-the-art results on both WN18RR and FB15k-237. On WN18RR, GMT achieves state-of-the-art performance with **0.621** MRR and **0.703** Hits@10, outperforming the strongest LLM-based baseline (MRR 0.593). On FB15k-237, GMT also ranks first with **0.488** MRR and **0.629** Hits@10, surpassing the best competitor (MRR 0.469). These con-

Table 4: The ablation results for the link prediction task.

Model	MRR	Hits@1	Hits@3	Hits@10
WN18RR				
GMT	0.621	0.569	0.667	0.703
w/o Semantics	0.603	0.541	0.639	0.695
w/o Fusion	0.558	0.516	0.607	0.639
FB15k-237				
GMT	0.488	0.394	0.505	0.629
w/o Semantics	0.464	0.368	0.479	0.602
w/o Fusion	0.425	0.331	0.450	0.584

sistent gains across both datasets validate the effectiveness of deep, memory-based graph evidence retrieval.

Triple Classification. We further evaluate GMT on triple classification over UMLS, CoDeX-S, and FB15K-237N, using a task-specific instruction template \mathcal{I}_{GMT} that differs from link prediction (Appendix A). As shown in Table 3, GMT delivers consistent state-of-the-art performance across the three benchmarks. In particular, GMT achieves the best

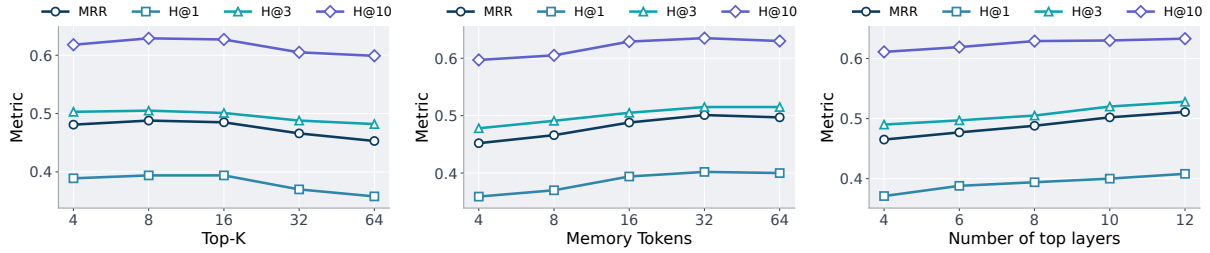


Figure 3: Impact of Key Hyperparameters on GMT Performance

overall results on UMLS (Acc **94.55**, F1 **93.76**), ranks first on CoDeX-S (Acc **89.01**, F1 **87.61**), and also leads on the more challenging FB15K-237N (Acc **84.10**, F1 **87.02**). These results demonstrate that GMT generalizes beyond link prediction and remains robust across datasets with diverse relational patterns and label distributions.

4.3 Ablation study: RQ2

To assess the contribution of each component, we compare GMT with two variants: **w/o Semantics**, which replaces SGM with conventional pre-trained KGE embeddings, and **w/o Fusion**, which removes the Graph-as-Memory Cross-Attention Fusion Module and injects memory by simple prefix concatenation.

Table 4 shows that both components are necessary. Removing SGM (**w/o Semantics**) consistently degrades performance on WN18RR and FB15k-237, indicating that relation-centric, context-aware semantics provide a stronger query-specific signal than static KG embeddings. More importantly, replacing cross-attention fusion with prefix injection (**w/o Fusion**) causes the largest drop (e.g., FB15k-237 MRR: 0.488 to 0.425), confirming that deep, token-wise memory retrieval is crucial and cannot be substituted by shallow concatenation.

Scoring Function in Stage-1 Pre-training. We further study the scoring function $F_G(h, r, t)$ used to pre-train SGM. As reported in Table 5, a RotatE-style scoring function performs best on FB15k-237 (MRR **0.471**), while alternatives such as TransE/DistMult/Complex or an MLP lead to lower accuracy. We therefore adopt RotatE by default in our implementation. The superiority of RotatE can likely be attributed to its inherent ability to model complex relational patterns.

Hyperparameter Sensitivity. We analyze the sensitivity of GMT to three key hyperparameters

Table 5: The scoring function for the link prediction task.

Model	MRR	Hits@1	Hits@3	Hits@10
FB15k-237				
GMT(w/ RotaE)	0.471	0.380	0.488	0.623
w TransE	0.459	0.375	0.483	0.610
w DistMult	0.454	0.366	0.481	0.615
w ComplexE	0.455	0.368	0.481	0.613
w MLP	0.447	0.358	0.475	0.606

Table 6: Impact of LLM Relational Knowledge Enhancement.

Model	MRR	Hits@1	Hits@3	Hits@10
FB15k-237				
GMT	0.488	0.394	0.505	0.629
w/o Enhancement	0.454	0.370	0.475	0.601

in Figure 3: (i) Top- K neighbor selection in SGM, (ii) the number of graph memory tokens m , and (iii) the number of top Transformer layers equipped with memory cross-attention. **Top- K** shows a clear sweet spot at a moderate neighborhood size (around $K=8-16$); too small loses context, while too large introduces noise and degrades MRR/Hits@1. Increasing the number of **memory tokens** improves performance and saturates around $m=32$, after which gains become marginal. Finally, injecting memory into more **top Transformer layers** consistently boosts all metrics, with the best results at the largest setting, supporting the benefit of deep, multi-layer memory retrieval.

4.4 Analysis of Semantics: RQ3

A key design in SGM is relation knowledge enhancement: we use an LLM to generate relation definitions and embed them as semantic vectors for Top- K neighbor selection. We evaluate **GMT w/o Enhancement**, which encodes only raw relation names with the same Sentence-BERT.

GMT w/o Enhancement. As shown in Table 6, removing enhancement consistently degrades performance on FB15k-237 (MRR drops from 0.488

Table 7: Case study: comparison of the top-5 neighbor relations for a query, before and after applying Knowledge Enhancement. Abbreviations stand for full relation names, e.g., **Gov Position (Title)** for `.../government_position_held/basic_title`. The Change column highlights the significant re-ranking based on semantic understanding.

Without Knowledge Enhancement (Lexical Matching)		With Knowledge Enhancement (Semantic Relevance)		
Neighbor Relation	Score	Neighbor Relation	Score	Change
Gov Position (Title)	0.859	Gov Position (Title)	0.611	—
Gov Position (Sessions)	0.828	Gov Position (Sessions)	0.460	—
Person (Profession)	0.250	Person (Nationality)	0.385	↑ Up (New)
Person (Places Lived)	0.206	Person (Employment)	0.377	↑ Up
Person (Employment)	0.184	Person (Profession)	0.338	↓ Down

Table 8: Robustness analysis across different LLM generators. Best results are in **bold**, second best are underlined.

Source	LLM Generator	WN18RR				FB15k-237				FB15k-237N			
		MRR	H@1	H@3	H@10	MRR	H@1	H@3	H@10	Acc	P	R	F1
Closed	GPT-4o (Ours)	0.621	<u>0.569</u>	0.667	0.703	<u>0.488</u>	0.394	<u>0.505</u>	<u>0.629</u>	84.10	83.14	<u>91.27</u>	87.02
	Claude-3.5-Sonnet	<u>0.617</u>	0.572	<u>0.660</u>	<u>0.702</u>	0.490	<u>0.390</u>	0.511	0.635	<u>84.05</u>	<u>83.11</u>	91.35	<u>86.98</u>
	Gemini-1.5-Pro	0.612	0.565	<u>0.656</u>	0.697	0.485	0.390	0.502	0.626	<u>83.95</u>	<u>82.98</u>	<u>91.10</u>	<u>86.85</u>
Open	Qwen3-32B-Instruct	0.619	0.567	0.664	0.700	0.486	0.391	0.503	0.627	83.92	82.95	91.08	86.82
	Qwen3-8B-Instruct	0.613	0.561	0.658	0.694	0.481	0.386	0.497	0.620	83.52	82.55	90.68	86.42
	Llama-3-8B-Instruct	0.610	0.558	0.655	0.690	0.478	0.383	0.494	0.616	83.25	82.26	90.35	86.11

to 0.454; Hits@1 drops from 0.394 to 0.370), verifying that relation names alone provide insufficient semantics for reliable neighbor filtering.

Robustness Across LLM Generators. We further investigate whether GMT relies on specific proprietary models for definition generation. Table 8 compares performance across various closed-source and open-source LLMs. Results show that GMT maintains consistent performance, with even lightweight open-source models yielding negligible drops. This confirms that GMT benefits from the explicit semantic guidance rather than the capabilities of a specific model.

Case Study. To highlight the role of relational semantics, we examine the query (Barack Obama, `/government/politician/government_positions_held...`, ?), which describes the geographical or administrative scope of a political position. As shown in Table 7, knowledge enhancement substantially re-orders the top-5 neighbors, shifting retrieval from surface lexical matching to semantic relevance. Without enhancement, the ranking is dominated by relations with shared string prefixes (e.g., Gov Position (Title) and Gov Position (Sessions)), which captures naming overlap rather than the intended semantics. With knowledge enhancement, neighbor retrieval becomes semantics-driven. In particular, Person (Nationality) appears and ranks highly, aligning

with the query’s notion of jurisdiction rather than string overlap. Meanwhile, semantically related relations (e.g., Person (Employment)) are promoted and less relevant ones (e.g., Person (Places Lived)) are suppressed, yielding a cleaner evidence set. Overall, the case study shows that SGM goes beyond lexical matching and builds a more semantically coherent graph memory, providing more reliable signals for LLM reasoning.

5 Conclusion

In this paper, we propose Graph-as-Memory Tuning (GMT), a memory-centric framework that integrates knowledge graphs with LLMs beyond shallow prefix fusion. GMT builds query-specific graph memory tokens using a Semantic Graph Module (SGM) with knowledge-enhanced relation semantics, and injects them into multiple Transformer layers via a Graph-as-Memory Cross-Attention Fusion Module, enabling token-wise retrieval of graph evidence during generation. We further adopt LoRA on memory cross-attention for parameter-efficient adaptation with a frozen base LLM.

Experiments on link prediction and triple classification benchmarks show that GMT achieves state-of-the-art or highly competitive performance, and ablations confirm the effectiveness of both SGM and cross-attention fusion, while additional analyses show GMT is robust to different LLMs used for relation knowledge enhancement.

540 Limitations

541 Despite its effectiveness, GMT still has sev-
542 eral method-specific limitations. First, its se-
543 mantic graph construction relies on an offline
544 knowledge-enhancement step that uses exter-
545 nal LLM-generated relation definitions and pre-
546 computed sentence embeddings, which introduces
547 external dependency and may make the relation se-
548 mantics less adaptive to highly context-dependent
549 queries.

550 Second, GMT builds graph memory from query-
551 local neighborhoods with Top-K semantic filtering,
552 so it may overlook useful long-range evidence, and
553 its performance remains sensitive to the neighbor-
554 hood size and memory configuration.

555 Third, GMT compresses variable-size subgraph
556 evidence into a fixed number of memory tokens and
557 aligns it with a frozen base LLM through LoRA
558 only on cross-attention, which may limit its ability
559 to preserve fine-grained graph information.

560 In future work, we plan to develop a more adap-
561 tive graph memory mechanism beyond the current
562 query-local neighborhood construction and fixed-
563 size memory tokenization, so that GMT can capture
564 longer-range structural dependencies and support
565 more complex multi-hop reasoning.

566 References

567 Ralph Abboud, Ismail Ilkan Ceylan, Thomas
568 Lukasiewicz, and Tommaso Salvatori. 2020. BoxE:
569 A box embedding model for knowledge base
570 completion. In *Proceedings of the 34th International
571 Conference on Neural Information Processing
572 Systems (NeurIPS 2020)*. Curran Associates Inc.

573 Antoine Bordes, Nicolas Usunier, Alberto Garcia-
574 Duran, Jason Weston, and Oksana Yakhnenko.
575 2013. Translating embeddings for modeling multi-
576 relational data. *Advances in neural information pro-
577 cessing systems*, 26.

578 Zongsheng Cao, Jing Li, Zigan Wang, and Jinliang Li.
579 2024. [Diffusione: Reasoning on knowledge graphs
580 via diffusion-based graph neural networks](#). In *Pro-
581 ceedings of the 30th ACM SIGKDD Conference on
582 Knowledge Discovery and Data Mining, KDD '24*,
583 page 222–230, New York, NY, USA. Association for
584 Computing Machinery.

585 Heng Chang, Jiangnan Ye, Alejo Lopez-Avila, Jinhua
586 Du, and Jia Li. 2024. Path-based explanation for
587 knowledge graph completion. In *Proceedings of the
588 30th ACM SIGKDD Conference on Knowledge Dis-
589 covery and Data Mining*, pages 231–242.

Zefeng Chen, Wensheng Gan, Jiayang Wu, Kaixia Hu,
and Hong Lin. 2025. Data scarcity in recommen-
dation systems: A survey. *ACM Transactions on
Recommender Systems*, 3(3):1–31. 590
591
592
593

Zexin Chen, Wenjie Peng, Zixuan Zhang, and Laks VS
Lakshmanan. 2024. ARR: A continuous-time dy-
namic KG embedding model for asynchronous and
recurring relations. In *Proceedings of the 2024 Inter-
national Conference on Management of Data (SIG-
MOD/PODS)*, pages 1–26. 594
595
596
597
598
599

Ziqiang Cui, Yunpeng Weng, Xing Tang, Fuyuan Lyu,
Dugang Liu, Xiuqiang He, and Chen Ma. 2025. [Com-
prehending knowledge graphs with large language
models for recommender systems](#). In *Proceedings
of the 48th International ACM SIGIR Conference on
Research and Development in Information Retrieval*,
SIGIR '25, page 1229–1239, New York, NY, USA.
Association for Computing Machinery. 600
601
602
603
604
605
606
607

Tim Dettmers, Pasquale Minervini, Pontus Stenetorp,
and Sebastian Riedel. 2018. Convolutional 2d knowl-
edge graph embeddings. In *Proceedings of the AAAI
Conference on Artificial Intelligence*, volume 32,
pages 3814–3820. 608
609
610
611
612

Jingtao Guo, Chunxia Zhang, Lingxi Li, Xiaojun Xue,
and Zhendong Niu. 2024a. A unified joint approach
with topological context learning and rule augmen-
tation for knowledge graph completion. In *Findings of
the Association for Computational Linguistics: ACL
2024*, Bangkok, Thailand. Association for Computa-
tional Linguistics. 613
614
615
616
617
618
619

Lingbing Guo, Zhongpu Bo, Zhuo Chen, Yichi Zhang,
Jiaoyan Chen, Yarong Lan, Mengshu Sun, Zhiqiang
Zhang, Yangyifei Luo, Qian Li, Qiang Zhang, Wen
Zhang, and Huajun Chen. 2024b. [Mkgl: Mastery
of a three-word language](#). In *Advances in Neural
Information Processing Systems*, volume 37, pages
140509–140534. Curran Associates, Inc. 620
621
622
623
624
625
626

Yuxin Jiang, Zixuan Zhang, Zexin Chen, Yun Tong,
Zonghan Wu, Nitesh V Chawla, and Laks VS Lak-
shmanan. 2024. KG-FIT: A framework of few-shot
inductive reasoning on knowledge graphs. In *Pro-
ceedings of the ACM Web Conference 2024 (WWW)*,
pages 111–122. 627
628
629
630
631
632

Binhang Li, Yizhou Liu, Zixuan Zhang, Rui Chen,
and Laks VS Lakshmanan. 2023. CSProm-KG: A
cross-modal self-supervised prompting for knowl-
edge graph completion under cold-start setting. In
*Findings of the Association for Computational Lin-
guistics: EMNLP 2023*, pages 11355–11368. 633
634
635
636
637
638

Qian Li, Zhuo Chen, Cheng Ji, Shiqi Jiang, and Jianxin
Li. 2024. [Llm-based multi-level knowledge gener-
ation for few-shot knowledge graph completion](#). In
*Proceedings of the Thirty-Third International Joint
Conference on Artificial Intelligence, IJCAI '24*. 639
640
641
642
643

Siyuan Li, Ruitong Liu, Te Sun, Yan Wen, Ruihao Zhou,
Jingyi Kang, and Yunjia Wu. 2025a. Context-driven
knowledge graph completion with semantic-aware 644
645
646

647	relational message passing. In <i>International Conference on Neural Information Processing</i> , pages 567–580. Springer.	703
648		704
649		705
650	Siyuan Li, Ruitong Liu, Yan Wen, Te Sun, Andi Zhang, Yanbiao Ma, and Xiaoshuai Hao. 2025b. Flow-modulated scoring for semantic-aware knowledge graph completion. <i>arXiv preprint arXiv:2506.23137</i> .	706
651		707
652		708
653		709
654	Qika Lin, Tianzhe Zhao, Kai He, Zhen Peng, Fangzhi Xu, Ling Huang, Jingying Ma, and Mengling Feng. 2025. Self-supervised quantized representation for seamlessly integrating knowledge graphs with large language models . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 13587–13602, Vienna, Austria. Association for Computational Linguistics.	710
655		711
656		712
657		713
658		714
659		715
660		716
661		717
662		718
663	Yu Liu, Yanan Cao, Xixun Lin, Yanmin Shang, Shi Wang, and Shirui Pan. 2025. Enhancing large language model for knowledge graph completion via structure-aware alignment-tuning. In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 20981–20995.	719
664		720
665		721
666		722
667		723
668		724
669	Yuyin Lu, Hegang Chen, Yanghui Rao, Jianxing Yu, Wen Hua, and Qing Li. 2025. An efficient fuzzy system for complex query answering on knowledge graphs . <i>IEEE Transactions on Knowledge and Data Engineering</i> , 37(9):4962–4976.	725
670		726
671		727
672		728
673		729
674	Kangyang Luo, Yuzhuo Bai, Cheng Gao, Shuzheng Si, Zhu Liu, Yingli Shen, Zhitong Wang, Cunliang Kong, Wenhao Li, Yufei Huang, Ye Tian, Xuantang Xiong, Lei Han, and Maosong Sun. 2025. GLTW: Joint improved graph transformer and LLM via three-word language for knowledge graph completion . In <i>Findings of the Association for Computational Linguistics</i> , Vienna, Austria. Association for Computational Linguistics.	730
675		731
676		732
677		733
678		734
679		735
680		736
681		737
682		738
683	Xin Lv, Yankai Lin, Yixin Cao, Lei Hou, Juanzi Li, Zhiyuan Liu, Peng Li, and Jie Zhou. 2022. Do pre-trained models benefit knowledge graph completion? a reliable evaluation and a reasonable approach. In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 3570–3581. Association for Computational Linguistics.	739
684		740
685		741
686		742
687		743
688		744
689		745
690	Reham Omar, Ishika Dhall, Panos Kalnis, and Essam Mansour. 2023. A universal question-answering platform for knowledge graphs. <i>Proceedings of the ACM on Management of Data</i> , 1(1):1–25.	746
691		747
692		748
693		749
694	Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Jipu Wang, and Xindong Wu. 2024. Unifying large language models and knowledge graphs: A roadmap . <i>IEEE Transactions on Knowledge and Data Engineering</i> , 36(7):3580–3599.	750
695		751
696		752
697		753
698		754
699	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.	755
700		756
701		757
702		758
	Zhiqing Sun, Zhi-Hong Deng, Jian-Yun Nie, and Jian Tang. 2019. Rotate: Knowledge graph embedding by relational rotation in complex space . In <i>International Conference on Learning Representations</i> .	707
		708
		709
		710
	Kristina Toutanova and Danqi Chen. 2015. Observed versus latent features for knowledge base and text inference. In <i>Proceedings of the 3rd Workshop on Continuous Vector Space Models and their Compositionality</i> , pages 57–66, Beijing, China. Association for Computational Linguistics.	711
		712
		713
		714
		715
		716
	Théo Trouillon, Johannes Welbl, Sebastian Riedel, Éric Gaussier, and Guillaume Cevaert. 2016. Complex embeddings for simple link prediction. In <i>Proceedings of the 33rd International Conference on Machine Learning (ICML)</i> , pages 2071–2080.	717
		718
		719
		720
		721
	Shikhar Vashishth, Soumya Sanyal, Vikram Nitin, and Partha Talukdar. 2020. Composition-based multi-relational graph convolutional networks . In <i>International Conference on Learning Representations</i> .	722
		723
		724
		725
	Hongwei Wang, Hongyu Ren, and Jure Leskovec. 2021. Relational message passing for knowledge graph completion . In <i>Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining, KDD '21</i> , page 1697–1707, New York, NY, USA. Association for Computing Machinery.	726
		727
		728
		729
		730
		731
	Linyi Wang, Wen Zhao, Zhipeng Wei, and Jing Liu. 2022. SimKGC: Simple contrastive knowledge graph completion with pre-trained language models. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4281–4294.	732
		733
		734
		735
		736
		737
	Yanbin Wei, Qiushi Huang, Yu Zhang, and James Kwok. 2023. KICGPT: Large language model with knowledge in context for knowledge graph completion . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 8667–8683, Singapore. Association for Computational Linguistics.	738
		739
		740
		741
		742
		743
	Hongcai Xu, Junpeng Bao, and Wenbo Liu. 2023. Double-branch multi-attention based graph neural network for knowledge graph completion. In <i>Proceedings of the 61st annual meeting of the association for computational linguistics (volume 1: long papers)</i> , pages 15257–15271.	744
		745
		746
		747
		748
		749
	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 <i>International Conference on Learning Representations (ICLR)</i> .	750
		751
		752
		753
		754
	Liang Yao, Chengsheng Mao, and Yuan Luo. 2019. Kg-bert: Bert for knowledge graph completion . <i>ArXiv</i> , abs/1909.03193.	755
		756
		757

758 Yichi Zhang, Zhuo Chen, Lingbing Guo, Yajing Xu,
759 Wen Zhang, and Huajun Chen. 2024a. Making large
760 language models perform better in knowledge graph
761 completion. In *Proceedings of the 32nd ACM inter-*
762 *national conference on multimedia*, pages 233–242.

763 Zequn Zhang, Jian Cai, Yanzhi Zhang, and Jie Wang.
764 2020. Learning hierarchy-aware knowledge graph
765 embeddings for link prediction. In *Proceedings of*
766 *the Thirty-Fourth AAAI Conference on Artificial In-*
767 *telligence*, pages 3065–3072. AAAI Press.

768 Zixuan Zhang, Ying Shao, Dmitri Kalashnikov, Binhang
769 Li, Yizhou Liu, and Laks VS Lakshmanan. 2024b.
770 Adaprop: A plug-and-play approach for propagating
771 large-scale updates in knowledge graphs. *Proceed-*
772 *ings of the VLDB Endowment*, 17(5):1016–1029.

773 Jia Zhu and Pasquale De Meo. 2025. [Exploring large](#)
774 [language models for knowledge graph completion](#)
775 [with auto-prompting](#). In *2025 International Joint*
776 *Conference on Neural Networks (IJCNN)*, pages 1–8.

777 Zhaocheng Zhu, Zuo Bai Zhang, Louis-Pascal Xhon-
778 neux, and Jian Tang. 2021. Neural bellman-ford
779 networks: A general graph neural network frame-
780 work for link prediction. In *Advances in Neural*
781 *Information Processing Systems*, volume 34, pages
782 29476–29490.

Appendix

A Prompt Template

Prompt: Knowledge Enhancement

Task: Generate a canonical and descriptive definition for a knowledge graph relation by following the provided examples. The definition must clearly explain the semantic relationship between a head entity and a tail entity in a neutral, factual sentence.

Relation Name: located_in

Definition: The head entity is a smaller geographical, physical, or conceptual area that is situated within the larger area of the tail entity.

Relation Name: relation_name

Definition: output

Prompt: Instruction Template for Link Prediction

Instruction: This is a knowledge graph completion task. Your goal is to predict the tail entity for an incomplete query triplet.

Problem Definition:

- **Query:** (head_entity, relation, ?)
- **Candidate Entities:** {candidate_1, ..., candidate_20}

Response Specification: Analyze the provided candidates and return a ranked list of the top three most likely answers, from most to least probable.

Prompt: Instruction Template for Triple Classification

Instruction: This is a triple classification task in knowledge graph. Your goal is to determine the correctness of the given triple.

Input:

- The triple is: (head_entity, relation, tail_entity)

Output: Please response True or False.