

KG-Attention: Knowledge Graph-Guided Attention at Test-Time via Bidirectional Information Aggregation

Songlin Zhai, Guilin Qi, Yuan Meng

School of Computer Science and Engineering, Southeast University, Nanjing, China
{songlin_zhai, gqi, yuan_meng}@seu.edu.cn

Abstract

Knowledge graphs (KGs) play a critical role in enhancing large language models (LLMs) by introducing structured and grounded knowledge into the learning process. However, most existing KG-enhanced approaches rely on parameter-intensive fine-tuning, which risks catastrophic forgetting and degrades the pretrained model’s generalization. Moreover, they exhibit limited adaptability to real-time knowledge updates due to their static integration frameworks. To address these issues, we introduce the first test-time KG-augmented framework for LLMs, built around a dedicated knowledge graph-guided attention (KGA) module that enables dynamic knowledge fusion without any parameter updates. The proposed KGA module augments the standard self-attention mechanism with two synergistic pathways: outward and inward aggregation. Specifically, the outward pathway dynamically integrates external knowledge into input representations via input-driven KG fusion. This inward aggregation complements the outward pathway by refining input representations through KG-guided filtering, suppressing task-irrelevant signals and amplifying knowledge-relevant patterns. Importantly, while the outward pathway handles knowledge fusion, the inward path selects the most relevant triples and feeds them back into the fusion process, forming a closed-loop enhancement mechanism. By synergistically combining these two pathways, the proposed method supports real-time knowledge fusion exclusively at test-time, without any parameter modification. Extensive experiments on five benchmarks verify the comparable knowledge fusion performance of KGA.

1 Introduction

Knowledge graphs (KGs) provide structured symbolic knowledge that complements the implicit parametric knowledge encoded in large language models (LLMs), giving rise to the emerging field of KG-enhanced LLMs [1, 2, 3, 4, 5]. While existing methods demonstrate promise in grounding LLMs with structured facts, most of them predominantly rely on parameter-invasive strategies such as various supervised fine-tuning strategies [1, 2, 3, 4, 6]. This fundamentally contradicts the preservation principle of LLMs [7]: parameter modification inevitably degrades acquired capabilities while introducing catastrophic forgetting. Furthermore, they struggle to adapt to real-time updates in knowledge graphs. An alternative paradigm is retrieval-augmented generation (RAG) [8, 9], which sidestep parameter updates but introduces new bottlenecks of retrieval reliability and latency [10]. Recent advancements in long-context LLMs enable full-context knowledge injection [11, 12, 13], yet they incur prohibitive time/memory overhead as context length increases, and also suffer from unpredictable interference from inter-triple dependencies [14]. These limitations motivate our central research question: *can we develop a parameter-preserving knowledge fusion method that achieves both efficient and real-time KG integration while maintaining LLMs’ general capabilities?*

Inspired by recent breakthroughs in *Test-Time Scalling* for LLMs [15, 16, 17], this paper introduces the knowledge graph-guided attention (**KGA**), a novel test-time adaptation framework that dynamically integrates external knowledge through non-invasive rewiring of attention interactions. Our design stems from a fundamental rethinking of the transformer attention, *i.e.*, the self-attention mechanism inherently supports adaptive information aggregation through token-to-token interactions [18, 19]. We extend this mechanism by reformulating “*self*”-aggregation in “*self*”-attention into a bidirectional “*knowledge*”-“*text*” synergistic aggregation in KGA via two well-designed pathways: (1) **Input→KG Interaction Flow** (Outward Aggregation) projects input queries to fuse KG triples, dynamically grounding predictions in external knowledge; and (2) **KG→Input Interaction Flow** (Inward Aggregation) recalibrates input representations to perform triple selection using KG-guided attention, suppressing noise while amplifying knowledge-critical patterns.

Specifically, the outward flow performs the interactions between input query representations and triple key/value features, thereby injecting external knowledge into input representations. The key-value matrices are transformed through parameter-shared attention module to ensure the alignment of feature space. Conversely, the inward flow applies triple-derived queries to aggregate input key/value representations, facilitating the triple selection in the outward aggregation. Crucially, both flows reuse the LLM’s native attention weights, preserving architectural integrity while enabling real-time knowledge updates via simple modifications to the triple strings. In addition, KGA also retains the original input→input self-attention flow for the preservation of linguistic and semantic understanding of inputs. We conduct extensive experiments on 5 benchmarks covering 3 different tasks using exclusively inference-time fusion, demonstrating the effectiveness of the proposed framework. To summarize, the contributions of this paper could be listed as:

- **Test-Time Knowledge Fusion:** A parameter-free mechanism enabling dynamic KG integration during inference via non-invasive attention-level rewiring.
- **Bidirectional Information Aggregation:** A novel method supporting mutual querying between original input and external KG information, ensuring adaptive knowledge fusion.
- **Comprehensive Empirical Validation:** Extensive evaluation across diverse tasks, demonstrating superior accuracy and scalability over 5 datasets.

2 Related Work

2.1 Knowledge Enhancement for LLMs

Large language models exhibit strong generative and reasoning capabilities, but still struggle with factual consistency and structured reasoning. To address these limitations, researchers augment LLMs with structured knowledge graphs, either through training-time integration or inference-time augmentation. Training-time methods inject KG information during pretraining or fine-tuning, such as K-BERT [20], ERNIE [2] and CoLAKE [21] which inject structured knowledge through input augmentation and unified pretraining objectives, respectively. Inference-time methods augment LLMs with external knowledge during decoding. For example, such methods augment LLMs via: leveraging mechanisms such as joint reasoning over KG subgraphs [22], multi-stage retrieval and inference pipelines [23], or non-parametric memory fusion for language generation [10]. Recent work such as KELP [24] also explores latent semantic matching to improve path-level knowledge selection. However, existing inference-time approaches often rely on static retrieval pipelines and lack mechanisms for fine-grained control over knowledge selection and integration.

2.2 Structured Attention Mask

Structured attention plays a pivotal role in equipping LLMs with structural awareness and scalable reasoning capabilities. Some methods focus on integrating external knowledge through attention masking. For example, KBLAM [25] uses rectangular masks to integrate key-value knowledge directly into the attention layer, eliminating the need for retrieval. Kformer [26] injects knowledge into the feedforward layer, providing a simple and effective alternative to attention-based integration. Zhu et al. [27] enhance the self-attention mechanism by incorporating AMR graph structure, allowing the model to capture dependencies between non-adjacent concepts in semantic generation tasks. Recent studies, including parallel context windows [28], structured contextual cues [29], leverage attention masks to isolate contextual fragments, improving modularity and inference efficiency.

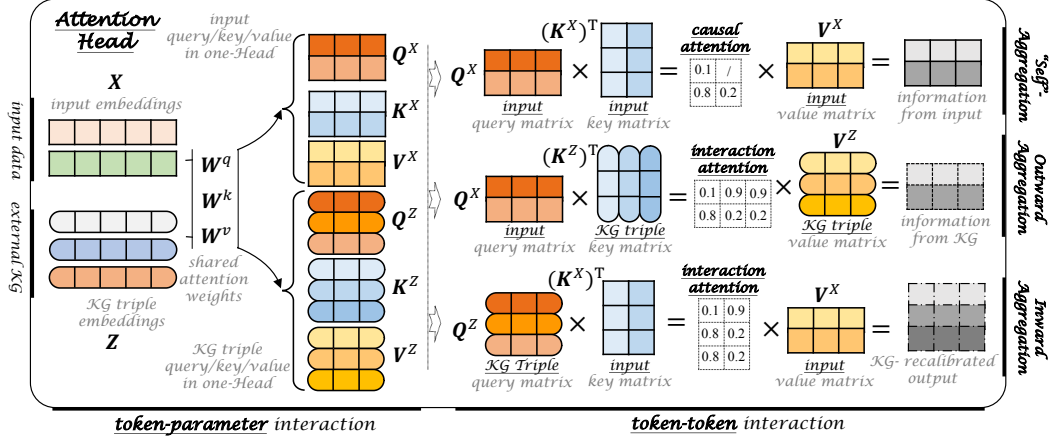


Figure 1: Illustration of the proposed knowledge graph-guided attention module.

3 Knowledge Graph-Guided Attention

3.1 Notations and Task Definition

Let \mathcal{M} denote a pre-trained large language model that receives a textual prompt X as input. The input X is tokenized into a sequence of tokens $X = \{x_1, \dots, x_n, \dots, x_N\}$ where N is the sequence length. At each transformer layer $l \in [1, L]$, the hidden representation of token x_n is denoted as the bold item $\mathbf{x}_n^{(l)}$, with $\mathbf{X}^{(l)} = [\mathbf{x}_1; \dots; \mathbf{x}_N]$ representing the full input feature matrix at layer l . For notational simplicity, the layer superscript is omitted (x_n, \mathbf{X}) when discussing operations within a single layer.

External knowledge \mathcal{G} is provided as a set of triples, we use $Z \in \mathcal{G}$ to denote a specific textual triple from a knowledge graph, e.g., (*NeurIPS_2025*, *held_in*, *San Diego*). Our goal is to develop an integration mechanism that could dynamically incorporate these triple knowledge into the forward computation of \mathcal{M} to enhance its ability to generate the desired output Y^* in response to the input X .

3.2 Bidirectional Information Aggregation: An Overview

To achieve the knowledge fusion goal, our primary contribution lies in the introduction of *knowledge graph-guided attention* module, which rewires the original self-attention module via establishing two extra synergistically enhanced information aggregation flows between input X and external knowledge \mathcal{G} , while preserving \mathcal{M} 's native parameters. Specifically, these two flows include: (1) **Outward Aggregation Pathway** (Input \rightarrow KG) dynamically injects external information into the input representations via the attention querying from X to \mathcal{G} . This pathway is termed “outward” as it focuses on the fusion of “external” knowledge. (2) **Inward Aggregation Pathway** (KG \rightarrow Input) strategically recalibrates input representations using the semantics of external knowledge, amplifying knowledge-critical patterns while suppressing irrelevant features. This pathway is designated “inward” due to its emphasis on “input” reorganization guided by external knowledge. Notably, the arrow notation \rightarrow specifically denotes the aggregation directionality, indicating the process of using source representations to reorganize target features, rather than data flow dependencies. Crucially, both interactions operate through parameter-preserving key-value transformations, avoiding any architectural modifications or new learnable parameters. Additionally, the framework retains the original “self”-attention flow (input \rightarrow input) as the foundational pathway to preserve the model’s intrinsic language modeling capabilities of X , forming a cohesive tri-flow architecture (see Figure 1).

3.3 Input→KG (Outward) Aggregation: External Knowledge Fusion

This pathway’s primary function is to dynamically integrate external information through input-driven knowledge aggregation. For a given triple Z , we first transform its string-form texts¹ into continuous D -dimensional hidden features through the embedding layer, generating *knowledge representations*. These representations are subsequently fused into input features via our KGA layer by layer. Formally, we demonstrate the fusion process using the l -th layer’s KGA (layer index l omitted in the feature matrix for clarity):

$$Q^Z = ZW^Q, \quad K^Z = ZW^K, \quad V^Z = ZW^V \quad (1)$$

where Z denotes the feature matrix of external knowledge triple Z containing sequential knowledge token representation $[z_1, \dots, z_m, \dots, z_M]$, with M representing the token count in triple Z . The derived components $q_m^Z \in Q^Z$, $k_m^Z \in K^Z$ and $v_m^Z \in V^Z$ respectively denote the query, key and value features of the m -th triple token. Crucially, W^Q , W^K and W^V are the original weight matrices in the self-attention module².

After obtaining the triple representations, the **outward aggregation** flow is then formulated as:

$$\hat{x}_n = \left\{ \overbrace{\sum_{i=1}^n \text{Exp}(w_{n,i}^X) v_i^X}^{\text{aggregation on inputs: } X} + \overbrace{\sum_{j=1}^M \text{Exp}\left(\frac{(q_n^X)^\top k_j^Z}{\sqrt{D}}\right) v_j^Z}^{\text{aggregation on external } Z \text{ via } X} \right\} / \left\{ \overbrace{\sum_{i=1}^n \text{Exp}(w_{n,i}^X) + \sum_{j=1}^M \text{Exp}\left(\frac{(q_n^X)^\top k_j^Z}{\sqrt{D}}\right)}^{\text{normalization}} \right\} \quad (2)$$

where $w_{n,i}^X = \{(q_n^X)^\top k_i^X\} / \sqrt{D}$ is the original attention score of input to itself. $q_n^X \in Q^X$ (distinct from the triple query feature q_m^Z) denotes the query vector of the n -th input token, with $k_n^X \in K^X$ and $v_n^X \in V^X$ representing the corresponding key and value vectors.

This outward aggregation enables input representations to actively “aggregate” information from external knowledge by the simple attention-level interaction, without any parameter updates. The pathway brings three key advantages: contextual adaptability (knowledge influence dynamically adjusts to input semantics), architectural purity (no interference with existing model architecture), and computational efficiency (linear scaling with triple length M). This fusion operates systematically across all transformer layers, exhibiting hierarchical specialization, *i.e.*, lexical-level fusion in lower layers, and semantic-level integration in higher layers.

3.4 KG→Input (Inward) Aggregation: Structured Input Recalibration

Building on the previous aggregation, we could achieve effective knowledge fusion. However, a key challenge lies in selecting relevant triples from \mathcal{G} to avoid exhaustive computation over the entire knowledge graph. To address this, we introduce a complementary aggregation flow. This pathway complements the outward aggregation by first carrying out a knowledge graph-guided refinement of input representations. This could strategically suppress irrelevant signals with amplifying knowledge-critical patterns in the input sequence, effectively addressing the potential attention dispersion noise amplification issue in conventional self-attention [31, 32]. The core interaction is formally defined as:

$$r_m = \left\{ \overbrace{\sum_{i=1}^N \text{Exp}\left(\frac{(q_m^Z)^\top k_i^X}{\sqrt{D}}\right) v_i^X}^{\text{aggregation on input } X \text{ by } Z} \right\} / \left\{ \overbrace{\sum_{i=1}^N \text{Exp}\left(\frac{(q_m^Z)^\top k_i^X}{\sqrt{D}}\right)}^{\text{normalization}} \right\} \quad (3)$$

where r_m denotes the recalibrated representation vector obtained by aggregating input representations X under the guidance of the m -th KG triple token z_m . This vector represents the remaining information after carefully filtering the input information using z_m . Crucially, this non-causal attention mechanism allows full input context visibility (unrestricted by causal masking), ensuring comprehensive contextual utilization during recalibration. The restructured input matrix $R = \{r_1, \dots, r_m, \dots, r_M\}$ encodes KG-refined input representations across all M triple tokens.

Adaptive Consolidation After obtaining the recalibrated token-level representations, we can sum these feature vectors (*i.e.*, $\{r_1 \dots r_m\}$) to obtain the input-level feature which is an aggregated input

¹We directly input the textual triple (*e.g.*, (*NeurIPS_2025*, *held_in*, *San Diego*)) to avoid error introduction from triple-to-text generation, due to their demonstrated capability of understanding triplet knowledge [14, 30].

²Direct parameter reuse ensures feature space alignment and architectural preservation.

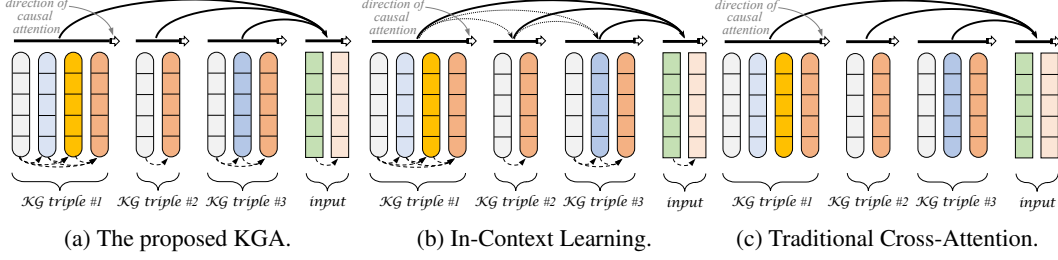


Figure 2: Illustrative depiction of information aggregation under different methods.

representation based on the entire triple. However, the importance of each token in the triple is not the same, which means that the importance of r_m should not be treated as the same. If a uniform summation is applied, the importance of individual components may be diluted, potentially obscuring critical information. As such, we propose an attention-aware fusion strategy:

$$\hat{r} = w_M^Z R = \sum_{j=1}^M w_{M,j}^Z r_j \quad (4)$$

Here, $w_{M,j}^Z$ represents the causal attention score between the M -th triple token (the last one) and the j -th token in Z , derived from the KG triple’s native self-attention. This design leverages the LLMs’ inherent knowledge of the triple token importance, while preserving architectural consistency.

Triple Selection The consolidated input representation \hat{r} derived from Eq. 4 enables systematic evaluation of input triple importance, achieving a self-refining feedback mechanism. This process guides knowledge triple selection in the outward aggregation phase, optimizing the overall knowledge fusion dynamics. Formally, we compute relevance between the triple and input as:

$$s(Z, X) = \hat{r}^\top x_N \quad (5)$$

This similarity score quantifies the triple’s contextual relevance to the input, enabling dynamic triple selection via top-k ranking. This gating mechanism ensures knowledge-critical filtering, *i.e.*, only triples with high contextual alignment influence final predictions.

3.5 Discussion

Section 3.3 delineates the external knowledge enhancement mechanism through input-driven triple integration (*i.e.*, outward aggregation). While this pathway shares superficial similarities with several conventional approaches, our framework exhibits fundamental architectural innovations. As shown in Figure 2a, the proposed outward aggregation systematically integrates three essential fusion dimensions: (1) Triple-to-Input: Knowledge infusion from KG to input sequence, (2) Triple-Itself: Intra-token self-attention for semantic consolidation in a given triple (3) Input-Itself: Native self-attention for contextual coherence maintenance. This tripartite architecture enables simultaneous refinement of both input representations and knowledge embeddings through their intrinsic patterns.

However, **In-Context Learning** in Figure 2b introduces uncontrolled inter-triple interference due to concatenating all triples as extended input. This forced co-processing of potentially unrelated triples (a) induces interference of different triples, (b) introduces sensitivity to triple ordering, and (c) unnecessary computation of inter-triple interaction and decreases the model efficiency.

Cross-Attention in Figure 2c restricts to unidirectional knowledge fusion (Triple→Input) while failing to preserve essential attention flows of the input itself (*i.e.*, the first term in Eq. 2). This architectural deficiency prevents effective knowledge refinement (no Triple-Itself) and degrades input contextualization (without Input-Itself).

4 Experiments

4.1 Experimental Settings

To validate the effectiveness of our knowledge fusion framework, we conduct experiments across three distinct tasks: single&multi-hop knowledge graph question answering (KGQA), knowledge

Table 1: Statistics of datasets.

Key	MetaQA			SimpleQuestions	PathQuestion 2-Hop	ZSRE EDIT	COUNTERFACT EDIT
	1-Hop	2-Hop	3-Hop				
#Test	9,947	14,872	14,274	21,687	1,908	19,085	10,000
Max($ \mathcal{G} $)/Case	380	11,850	89,522	620,668	188	/	/
Min($ \mathcal{G} $)/Case	1	2	7	1	2	/	/

graph reasoning, and knowledge-based model editing (KME). For KGQA and KG reasoning tasks, models are required to integrate knowledge triples from external KGs into input (see Figure 1). While naively feeding the entire KG to the model provides the most straightforward solution, this incurs prohibitive computational overhead and is unnecessary. Instead, we first perform entity linking via the ordinary string similarity matching³ on question mentions to avoid dependency on external modules, then retrieve all related triples using the identified *Top-1* entity as an anchor. For instance, in the **Single-hop KGQA**, we retrieve all triples directly connected to the anchor entity as candidate triples. In the **Multi-hop KGQA**, we extend retrieval to 2-hop neighbors of the anchor entity. Table 1 summarizes the statistics of candidate triples, including minimum and maximum counts per instance. For knowledge-based model editing task (KME), where target knowledge modifications is predefined, our method directly fuses provided triples without additional screening. Notably, we consider a more practical yet challenging scenario, *consecutive editing*, involving sequential knowledge updates without parameter rollback after each edit. Specifically, we evaluate model performance under 3,000 edits to assess long-term stability. All experiments are conducted based on the Qwen2.5 7B model.

For comparative evaluation in KGQA and KG Reasoning, we compare against three categories of baselines: (1) **Supervised Fine-Tuned Methods (SFT)**: representing traditional supervised approaches with KG-augmented training, (2) **In-Context Learning (ICL)**⁴: concatenating all candidate triples as contextual prefixes to input questions⁵, (3) **Zero-Shot Learning (ZSL)**: directly processing input questions without KG integration to establish reference performance levels. This experimental setup ensures rigorous validation of our method’s ability to dynamically integrate knowledge while preserving architectural integrity, with all baselines implemented under identical hardware (NVIDIA A100 80G) and evaluation protocols.

4.2 Overall Performance Comparison

Table 2 provides a comprehensive comparison of KGA’s performance across different datasets, revealing three critical insights. First, SFT-based models gain satisfying results, even achieving near-perfect results on KGQA (*e.g.*, PullNet attains >99% accuracy on MetaQA 2-Hop). While this type of models fit the current data well, extensive studies reveal that they still suffer from severe limitations: **poor cross-domain transferability** [5] (*e.g.*, 40+ % accuracy drop in Wikidata→DBpedia transfers), **weak adaptability to KG dynamics** [68], and **inability to support incremental learning**. These limitations restrict the practical applicability of such models. Second, while ZSL yields poor performance (*e.g.*, 19.7% on MetaQA 1-Hop), ICL significantly improves accuracy (*e.g.*, 71.1% on MetaQA 1-Hop) but introduces prohibitive computational costs and memory overhead (detailed in Sections 4.4 and 4.5) due to unfiltered triple aggregation. A critical limitation stems from ICL’s unfiltered processing of all candidate triples as contextual input. By feeding the entire triple set into the self-attention mechanism without explicit screening mechanisms, the model inevitably aggregates information from irrelevant triples, *e.g.*, spuriously linking “(Donald Trump, held_position, 47 POTUS)→(NeurIPS_2025, held_in, San Diego)”, introducing two key issues: **suboptimal computational efficiency** due to redundant attention computations over non-critical triples, and **degraded interpretability** from cross-triple token interactions that obscure actual knowledge utilization patterns [14]. These drawbacks underscore the necessity of our calibrated aggregation approach for precise knowledge-text alignment. Third, our KGA framework demonstrates competitive advan-

³The core of our method lies in fusion rather than graph retrieval. This low-coupling design allows seamless integration with advanced triple screening mechanisms to enhance performance and efficiency.

⁴Follow the paradigm established in Wang et al. [14].

⁵Candidate triples are limited to 100 per instance for computational resource traceability, and experiments are conducted multiple times.

Table 2: Comparison between KGA and previous methods on different tasks, where methods with [†] are the ones based on GPT-3.5 and the ones with * are based on DeepSeek-V3. ZSL denotes the *Zero-Shot Learning* method that directly queries the LLM without providing any additional context, *i.e.*, merely using the LLM’s internal knowledge for answering. ICL indicates *In-Context Learning*, which flattens all KB triples into a natural language utterance, and attaches it in front of the prompt.

(Multi&Single -Hop Knowledge Graph Question Answering)					(KG Reasoning)		
MetaQA				SimpleQuestions		PathQuestion	
Methods (Hit@1)	1-Hop	2-Hop	3-Hop	Methods	Hit@1	Methods	Hit@1
Bordes [33]	95.7	81.8	28.4	MemNN [34]	63.9	ISM [35]	99.1
KV-MemNN [36]	95.8	25.1	10.1	CFO [37]	62.6	QAGCN [38]	98.5
VRN [39]	97.5	89.9	62.5	AMPCNN [40]	67.2	Uhop-HR [41]	97.6
EmbedKGQA [1]	97.5	98.8	94.8	Character [42]	70.3	KV-MemNN [36]	97.4
GraftNet [43]	97.0	94.8	77.7	IOPrompt [44]	20.0	SRN [45]	96.3
KG-GPT _{12shot} [6]	96.3	94.4	94.0	SC [46]	18.9	IRN [47]	96.0
PullNet [48]	97.0	99.9	91.4	RoG [49]	73.3	MINERVA [50]	75.9
SRN [45]	97.0	95.1	75.2	StructGPT [†] [51]	50.2	TransferNet [52]	93.2
G-Riever [53]	98.5	87.6	54.9	PoG* [54]	63.9	RL-MHR [55]	94.1
KAPING [56]	90.8	71.2	43.0	KnowPath* [57]	65.3	AlAgha [41]	97.4
SimGRAG [58]	98.0	98.4	97.8	ToG* [59]	59.7	ARN_DistMult	84.3
ZSL	19.7	20.7	22.6	ZSL	0.7	ZSL	25.1
ICL	71.1	49.4	27.9	ICL	4.2	ICL	73.7
KGA	80.9	68.3	43.2	KGA	5.8	KGA	80.3

(Knowledge-based Model Editing)

ZsRE					COUNTERFACT			
Methods	Efficacy	Generality	Locality	Score	Efficacy	Generality	Locality	Score
Qwen2.5 (7B)	20.2	19.4	/	19.8	1.3	0.4	/	0.9
FT-C [60]	2.4	2.0	0.3	1.6	4.1	1.5	0.1	1.9
LoRA [61]	5.3	5.2	0.7	3.7	2.0	2.0	0.2	1.4
ROME [62]	28.1	25.4	9.2	20.9	0.0	0.0	2.4	0.8
R-ROME [63]	59.1	50.2	38.7	49.3	57.8	17.4	45.9	40.4
MEMIT [64]	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
AlphaEdit [65]	0.01	0.01	0.02	0.01	0.0	0.0	0.0	0.0
IKE [66]	69.7	68.9	50.3	63.0	50.1	22.7	56.9	43.2
GRACE [67]	29.1	0.9	100.0	43.3	0.1	0.0	99.4	33.2
WISE [7]	27.5	26.9	9.9	21.4	3.5	2.8	6.4	4.2
KGA	75.9	73.8	100.0	83.2	69.3	50.9	100.0	73.4

tages: it outperforms ICL by 18.9% on MetaQA 2-Hop through active triple filtering, preserves full model capabilities (evidenced by 100% *Locality* in KME), and enables dynamic KG adaptation through simple updates to the input triples. The integrated inward aggregation mechanism provides quantifiable interpretability (see Section 4.6) by tracing layer-wise attention to critical triples, while avoiding unnecessary computations through hierarchical screening. These results position KGA as a balanced solution for practical deployments requiring both knowledge awareness and computational efficiency. Notably, both KGA and ICL exhibit unexpectedly poor performance on the SimpleQuestions benchmark. Our error analysis reveals that this primarily stems from error propagation caused by imperfect candidate triplet retrieval via lexical matching, a limitation of current retrieval systems rather than our core methodology. Crucially, our framework remains decoupled from retrieval filtering mechanisms, enabling seamless integration with advanced retrieval techniques like dense semantic search (*e.g.*, DPR [69]) or cross-encoder reranking to mitigate such issues.

4.3 Effects of Inward Aggregation

As detailed in Section 3.4, the inward aggregation module re-aggregates input information, thereby obtaining a calibrated representation to filter candidate triples retrieved through exhaustive graph

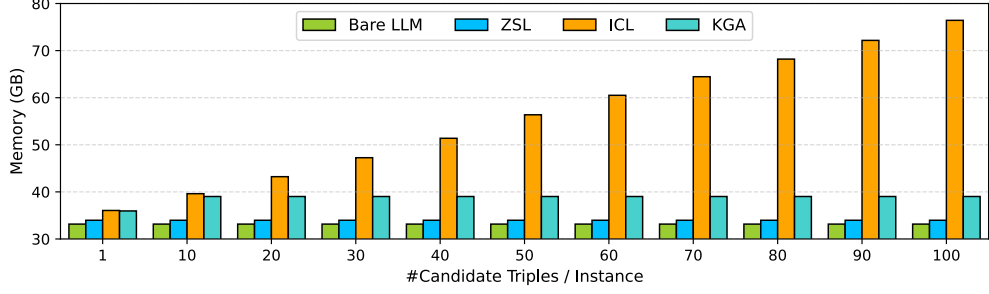


Figure 3: Comparison of memory usage under different methods.

traversal and enabling precise knowledge integration. To quantify its efficacy, we evaluate ground truth triple recall on the PathQuestion benchmark using:

$$\text{Recall}(\text{KGA}(\mathcal{G}, X)|X, \alpha) = \frac{1}{|\mathcal{X}|} \sum_{X \in \mathcal{X}} |\hat{\pi}(X, \alpha)| / |\pi(X)| \quad (6)$$

where $\pi(X)$ denotes the ground truth triple list for query X (with $|\pi(X)|$ termed unit length). $\hat{\pi}(X, \alpha)$ represents triples ranked within the Top- $\alpha|\pi(X)|$ positions by KGA. Figure 4 compares three strategies using this metric: (1) inward aggregation-enabled ranking, (2) module-disabled screening, and (3) random triple selection. It can be observed that this module achieves nearly 100% recall within 3 unit lengths, far smaller than the candidate triple pool size, significantly reducing first-stage fusion latency. In contrast, random selection attains $< 50\%$ recall even at $5 \times$ unit lengths, while the module-disabled variant underperforms by about 10% absolute recall, demonstrating the necessity of input recalibration for effective screening. These results validate inward aggregation’s dual role: eliminating computational redundancy while ensuring critical knowledge retention.

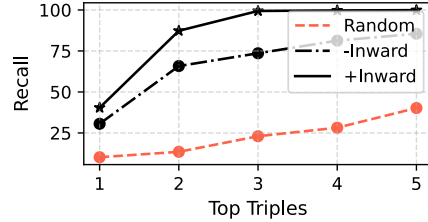


Figure 4: Illustration of the recall of three different triple selection strategies.

4.4 Efficiency Comparison

Our framework leverages external knowledge to recalibrate input representations for candidate triple screening, thereby enriching contextual relevance while eliminating redundant computations in the attention mechanism. To assess practical viability, we measure inference latency across three methods under varying candidate triple counts (Figure 5). ZSL baseline maintains stable latency (0.1 $s/instance$ on PathQuestion) as it processes no triples. In contrast, ICL exhibits quadratic time growth with triple count, *e.g.*, processing 100 triples incurs about $20 \times$ higher latency than ZSL due to indiscriminate attention over all input triples. However, some of these candidate triples are irrelevant to the current query, leading to wasted computation on unnecessary triple aggregation in the attention module. Our KGA framework circumvents this inefficiency through input-calibrated triple filtering: by fixing the processed triple count at an optimal threshold (30 for PathQuestion) once ground truth recall plateaus (Section 4.3), inference time stabilizes despite increasing candidate pools. While inward aggregation introduces preprocessing overhead proportional to triple count, this cost is amortizable through offline caching, reducing online latency to ZSL-comparable levels (0.73 $s/instance$). These results demonstrate KGA’s ability to balance computational efficiency and knowledge fidelity.

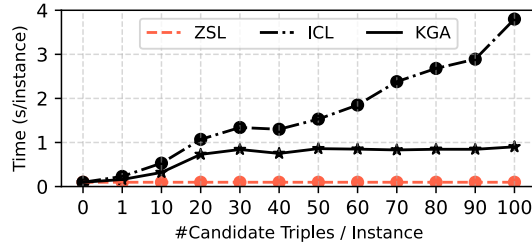


Figure 5: Efficiency of three different methods.

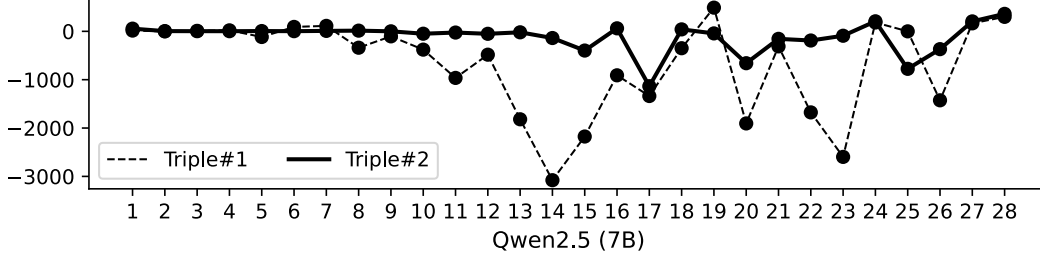


Figure 6: Dynamics of different triples’ importance across different hidden layers, where Triple#1 refers to (*NeurIPS_2025, held_in, San Diego*), and Triple#2 represents (*San Diego, located_in, USA*), with the query being *In which country is NeurIPS_2025 held?*.

4.5 Memory Usage

Memory footprint constitutes a critical deployment consideration for knowledge fusion systems. We evaluate GPU memory consumption across three methods under varying input triple counts, benchmarking against the base model’s requirement (32 GB). The ZSL baseline maintains stable memory usage (~ 33 GB) as it processes no triples, while ICL exhibits dramatic memory growth, expanding from 33 GB (1 triple) to 78 GB (100 triples). The memory requirement of ICL will further escalate for token-dense triples (*e.g.*, 112+ GB for 100 verbose triples). This overhead stems largely from caching non-critical triple features, as analyzed in Section 4.4. However, our method addresses this by filtering triples through recalibrated input representations: by fixing the processed triple count at an optimal threshold, we achieve predictable memory consumption independent of total candidate volume. This stability enables deployment on resource-constrained edge devices, positioning our framework as a practical solution for memory-sensitive knowledge applications.

4.6 Case Study

To clarify how our model utilizes input triples and enhance explainability, we analyze a concrete example: answering the question “*In which country is NeurIPS_2025 held?*” with two input triples (*NeurIPS_2025, held_in, San Diego*) and (*San Diego, located_in, USA*). This two-hop reasoning case requires synthesizing information from both triples. Figure 6 demonstrates the layer-wise attention weights allocated to each triple (as computed by Eq. 5). From this figure, we can observe that: **In shallow layers (1-4)**, both triples receive comparable attention as the model prioritizes linguistic feature extraction (*e.g.*, entity recognition and syntactic parsing). **Intermediate layers (4-7)** show slightly higher attention to the first triple due to its direct mention of the query entity *NeurIPS_2025*. **Deeper layers (8-17)** exhibit significantly increased focus on the second triple, reflecting its critical role in resolving the geographic relationship -level inference (*located_in* \rightarrow *country*). **Final layers (18-28)** display oscillating attention patterns, indicating iterative refinement and cross-triple evidence integration during answer generation. This hierarchical attention progression aligns with LLMs’ architectures’ characteristic processing: shallow feature grounding \rightarrow intermediate semantic association \rightarrow deep relational reasoning \rightarrow final prediction synthesis [70, 71].

5 Conclusion

This paper introduced a novel test-time knowledge fusion framework that dynamically integrates external knowledge graphs into large language models through parameter-preserving attention rewiring, eliminating the need for fine-tuning or architectural modifications. By establishing bidirectional interactions between textual inputs and triple tokens via outward and inward aggregation pathways, our method achieved competitive performance across knowledge-intensive tasks while maintaining the base model’s general capabilities and computational efficiency. Experimental validation demonstrated its effectiveness, positioning the framework as a practical solution for deploying reliable, knowledge-aware LLMs in real-world scenarios. The interpretable attention patterns and modular design further enable seamless integration with evolving KGs, bridging the gap between parametric knowledge and symbolic reasoning without compromising deployment scalability.

References

- [1] Apoorv Saxena, Aditay Tripathi, and Partha Talukdar. Improving multi-hop question answering over knowledge graphs using knowledge base embeddings. In *ACL*, pages 4498–4507. Association for Computational Linguistics, 2020. URL <https://aclanthology.org/2020.acl-main.412/>.
- [2] Yu Sun, Shuohuan Wang, Shikun Feng, Siyu Ding, Chao Pang, Junyuan Shang, Jiaxiang Liu, Xuyi Chen, Yanbin Zhao, Yuxiang Lu, Weixin Liu, Zihua Wu, Weibao Gong, Jianzhong Liang, Zhizhou Shang, Peng Sun, Wei Liu, Xuan Ouyang, Dianhai Yu, Hao Tian, Hua Wu, and Haifeng Wang. ERNIE 3.0: Large-scale knowledge enhanced pre-training for language understanding and generation. *CoRR*, abs/2107.02137, 2021. URL <https://arxiv.org/abs/2107.02137>.
- [3] Xiao Huang, Jingyuan Zhang, Dingcheng Li, and Ping Li. Knowledge graph embedding based question answering. In *WSDM*, page 105–113. Association for Computing Machinery, 2019. URL <https://doi.org/10.1145/3289600.3290956>.
- [4] Jinhao Jiang, Kun Zhou, Xin Zhao, and Ji-Rong Wen. Unikgqa: Unified retrieval and reasoning for solving multi-hop question answering over knowledge graph. In *ICLR*. OpenReview.net, 2023. URL <https://openreview.net/forum?id=Z63RvyAZ2Vh>.
- [5] Haoran Luo, Haihong E, Zichen Tang, Shiyao Peng, Yikai Guo, Wentai Zhang, Chenghao Ma, Guanting Dong, Meina Song, Wei Lin, Yifan Zhu, and Anh Tuan Luu. ChatKBQA: A generate-then-retrieve framework for knowledge base question answering with fine-tuned large language models. In *Findings of the ACL*, pages 2039–2056. Association for Computational Linguistics, 2024. URL <https://aclanthology.org/2024.findings-acl.122/>.
- [6] Jiho Kim, Yeonsu Kwon, Yohan Jo, and Edward Choi. KG-GPT: A general framework for reasoning on knowledge graphs using large language models. In *EMNLP*, pages 9410–9421. Association for Computational Linguistics, 2023. URL <https://aclanthology.org/2023.findings-emnlp.631/>.
- [7] Peng Wang, Zexi Li, Ningyu Zhang, Ziwen Xu, Yunzhi Yao, Yong Jiang, Pengjun Xie, Fei Huang, and Huajun Chen. WISE: Rethinking the knowledge memory for lifelong model editing of large language models. In *NeurIPS*, 2024. URL <https://openreview.net/forum?id=VJMY0fJVC2>.
- [8] Wenyu Huang, Guancheng Zhou, Hongru Wang, Pavlos Vougiouklis, Mirella Lapata, and Jeff Z. Pan. Less is more: Making smaller language models competent subgraph retrievers for multi-hop KGQA. In *Findings of EMNLP*, pages 15787–15803. Association for Computational Linguistics, 2024. URL <https://aclanthology.org/2024.findings-emnlp.927/>.
- [9] Diego Sanmartin. KG-RAG: bridging the gap between knowledge and creativity. *CoRR*, abs/2405.12035, 2024. URL <https://doi.org/10.48550/arXiv.2405.12035>.
- [10] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *NeurIPS*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html>.
- [11] Zhuowan Li, Cheng Li, Mingyang Zhang, Qiaozhu Mei, and Michael Bendersky. Retrieval augmented generation or long-context LLMs? a comprehensive study and hybrid approach. In *EMNLP Industry Track*. Association for Computational Linguistics, 2024. URL <https://aclanthology.org/2024.emnlp-industry.66/>.
- [12] Tan Yu, Anbang Xu, and Rama Akkiraju. In defense of RAG in the era of long-context language models. *CoRR*, abs/2409.01666, 2024. URL <https://doi.org/10.48550/arXiv.2409.01666>.
- [13] Gemini Team. Gemini: A family of highly capable multimodal models, 2024. URL <https://arxiv.org/abs/2312.11805>.

- [14] Xi Wang, Taketomo Isazawa, Liana Mikaelyan, and James Hensman. KBLam: Knowledge base augmented language model. In *ICLR*, 2025. URL <https://openreview.net/forum?id=aLsMzkTej9>.
- [15] Zhiyuan Zeng, Qinyuan Cheng, Zhangyue Yin, Yunhua Zhou, and Xipeng Qiu. Revisiting the test-time scaling of ol-like models: Do they truly possess test-time scaling capabilities?, 2025. URL <https://arxiv.org/abs/2502.12215>.
- [16] Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. sl: Simple test-time scaling, 2025. URL <https://arxiv.org/abs/2501.19393>.
- [17] William Jurayj, Jeffrey Cheng, and Benjamin Van Durme. Is that your final answer? test-time scaling improves selective question answering, 2025. URL <https://arxiv.org/abs/2502.13962>.
- [18] Jean-Baptiste Cordonnier, Andreas Loukas, and Martin Jaggi. On the relationship between self-attention and convolutional layers. In *ICLR*. OpenReview.net, 2020. URL <https://openreview.net/forum?id=HJlnC1rKPB>.
- [19] Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned. In *ACL*, pages 5797–5808. Association for Computational Linguistics, 2019. URL <https://doi.org/10.18653/v1/p19-1580>.
- [20] Weijie Liu, Peng Zhou, Zhe Zhao, Zhiruo Wang, Qi Ju, Haotang Deng, and Ping Wang. K-BERT: enabling language representation with knowledge graph. In *AAAI*, pages 2901–2908. AAAI Press, 2020. URL <https://doi.org/10.1609/aaai.v34i03.5681>.
- [21] Tianxiang Sun, Yunfan Shao, Xipeng Qiu, Qipeng Guo, Yaru Hu, Xuanjing Huang, and Zheng Zhang. Colake: Contextualized language and knowledge embedding. In *COLING*, pages 3660–3670. International Committee on Computational Linguistics, 2020. URL <https://doi.org/10.18653/v1/2020.coling-main.327>.
- [22] Michihiro Yasunaga, Hongyu Ren, Antoine Bosselut, Percy Liang, and Jure Leskovec. QA-GNN: reasoning with language models and knowledge graphs for question answering. In *NAACL-HLT*, pages 535–546. Association for Computational Linguistics, 2021. URL <https://doi.org/10.18653/v1/2021.naacl-main.45>.
- [23] Jiho Kim, Yeonsu Kwon, Yohan Jo, and Edward Choi. KG-GPT: A general framework for reasoning on knowledge graphs using large language models. In *Findings of EMNLP*, pages 9410–9421. Association for Computational Linguistics, 2023. URL <https://doi.org/10.18653/v1/2023.findings-emnlp.631>.
- [24] Haochen Liu, Song Wang, Yaochen Zhu, Yushun Dong, and Jundong Li. Knowledge graph-enhanced large language models via path selection. In *Findings of ACL*, pages 6311–6321. Association for Computational Linguistics, 2024. URL <https://doi.org/10.18653/v1/2024.findings-acl.376>.
- [25] Xi Wang, Liana Mikaelyan, Taketomo Isazawa, and James Hensman. Kblam: Knowledge base augmented language model. *CoRR*, abs/2410.10450, 2024. URL <https://doi.org/10.48550/arXiv.2410.10450>.
- [26] Yunzhi Yao, Shaohan Huang, Li Dong, Furu Wei, Huajun Chen, and Ningyu Zhang. Kformer: Knowledge injection in transformer feed-forward layers. In *NLPCC*, volume 13551 of *Lecture Notes in Computer Science*, pages 131–143. Springer, 2022. URL https://doi.org/10.1007/978-3-031-17120-8_11.
- [27] Jie Zhu, Junhui Li, Muhua Zhu, Longhua Qian, Min Zhang, and Guodong Zhou. Modeling graph structure in transformer for better amr-to-text generation. In *EMNLP-IJCNLP*, pages 5458–5467. Association for Computational Linguistics, 2019. URL <https://doi.org/10.18653/v1/D19-1548>.

- [28] Nir Ratner, Yoav Levine, Yonatan Belinkov, Ori Ram, Inbal Magar, Omri Abend, Ehud Karpas, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. Parallel context windows for large language models. In *ACL (Volume 1: Long Papers)*, pages 6383–6402. Association for Computational Linguistics, 2023.
- [29] Thomas Merth, Qichen Fu, Mohammad Rastegari, and Mahyar Najibi. Superposition prompting: Improving and accelerating retrieval-augmented generation. In *ICML*. OpenReview.net, 2024. URL <https://openreview.net/forum?id=r8k5JrGip6>.
- [30] Xinbang Dai, Yuncheng Hua, Tongtong Wu, Yang Sheng, Qiu Ji, and Guilin Qi. Large language models can better understand knowledge graphs than we thought. *Knowl. Based Syst.*, 312: 113060, 2025. URL <https://doi.org/10.1016/j.knosys.2025.113060>.
- [31] Murtadha Ahmed, Wenbo, and Liu yunfeng. Mateicl: Mitigating attention dispersion in large-scale in-context learning, 2025. URL <https://arxiv.org/abs/2505.01110>.
- [32] Tianzhu Ye, Li Dong, Yuqing Xia, Yutao Sun, Yi Zhu, Gao Huang, and Furu Wei. Differential transformer. In *ICLR*, 2025. URL <https://openreview.net/forum?id=0voCm1gGhN>.
- [33] Antoine Bordes, Sumit Chopra, and Jason Weston. Question answering with subgraph embeddings. In *EMNLP*, pages 615–620. Association for Computational Linguistics, 2014. URL <https://aclanthology.org/D14-1067/>.
- [34] Antoine Bordes, Nicolas Usunier, Sumit Chopra, and Jason Weston. Large-scale simple question answering with memory networks. *CoRR*, abs/1506.02075, 2015. URL <http://arxiv.org/abs/1506.02075>.
- [35] Yunshi Lan, Shuohang Wang, and Jing Jiang. Multi-hop knowledge base question answering with an iterative sequence matching model. In *ICDM*, pages 359–368. IEEE, 2019. URL <https://doi.org/10.1109/ICDM.2019.00046>.
- [36] Alexander Miller, Adam Fisch, Jesse Dodge, Amir-Hossein Karimi, Antoine Bordes, and Jason Weston. Key-value memory networks for directly reading documents. In *EMNLP*, pages 1400–1409. Association for Computational Linguistics, 2016. URL <https://aclanthology.org/D16-1147/>.
- [37] Zihang Dai, Lei Li, and Wei Xu. CFO: conditional focused neural question answering with large-scale knowledge bases. In *ACL*. The Association for Computer Linguistics, 2016. URL <https://doi.org/10.18653/v1/p16-1076>.
- [38] Ruijie Wang, Luca Rossetto, Michael Cochez, and Abraham Bernstein. QAGCN: answering multi-relation questions via single-step implicit reasoning over knowledge graphs. In *ESWC*, volume 14664 of *Lecture Notes in Computer Science*, pages 41–58. Springer, 2024. URL https://doi.org/10.1007/978-3-031-60626-7_3.
- [39] Yuyu Zhang, Hanjun Dai, Zornitsa Kozareva, Alexander J. Smola, and Le Song. Variational reasoning for question answering with knowledge graph. In *AAAI*, pages 6069–6076. AAAI Press, 2018. URL <https://doi.org/10.1609/aaai.v32i1.12057>.
- [40] Wenpeng Yin, Mo Yu, Bing Xiang, Bowen Zhou, and Hinrich Schütze. Simple question answering by attentive convolutional neural network. In *COLING*, pages 1746–1756. ACL, 2016. URL <https://aclanthology.org/C16-1164/>.
- [41] Zi-Yuan Chen, Chih-Hung Chang, Yi-Pei Chen, Jijnasa Nayak, and Lun-Wei Ku. UHop: An unrestricted-hop relation extraction framework for knowledge-based question answering. In *NAACL*, pages 345–356. Association for Computational Linguistics, 2019. URL <https://aclanthology.org/N19-1031/>.
- [42] Xiaodong He and David Golub. Character-level question answering with attention. In *EMNLP*, pages 1598–1607. The Association for Computational Linguistics, 2016. URL <https://doi.org/10.18653/v1/d16-1166>.

- [43] Haitian Sun, Bhuwan Dhingra, Manzil Zaheer, Kathryn Mazaitis, Ruslan Salakhutdinov, and William Cohen. Open domain question answering using early fusion of knowledge bases and text. In *EMNLP*, pages 4231–4242. Association for Computational Linguistics, 2018. URL <https://aclanthology.org/D18-1455/>.
- [44] 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. Language models are few-shot learners. In *NeurIPS*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfc4967418bfb8ac142f64a-Abstract.html>.
- [45] Yunqi Qiu, Yuanzhuo Wang, Xiaolong Jin, and Kun Zhang. Stepwise reasoning for multi-relation question answering over knowledge graph with weak supervision. In *WSDM*, page 474–482. Association for Computing Machinery, 2020. URL <https://doi.org/10.1145/3336191.3371812>.
- [46] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In *ICLR*. OpenReview.net, 2023. URL <https://openreview.net/forum?id=1PL1NIMMrw>.
- [47] Mantong Zhou, Minlie Huang, and Xiaoyan Zhu. An interpretable reasoning network for multi-relation question answering. In *COLING*, pages 2010–2022. Association for Computational Linguistics, 2018. URL <https://aclanthology.org/C18-1171/>.
- [48] Haitian Sun, Tania Bedrax-Weiss, and William Cohen. PullNet: Open domain question answering with iterative retrieval on knowledge bases and text. In *EMNLP-IJCNLP*, pages 2380–2390. Association for Computational Linguistics, 2019. URL <https://aclanthology.org/D19-1242/>.
- [49] Linhao Luo, Yuan-Fang Li, Gholamreza Haffari, and Shirui Pan. Reasoning on graphs: Faithful and interpretable large language model reasoning. In *ICLR*. OpenReview.net, 2024. URL <https://openreview.net/forum?id=ZGNWW7xZ6Q>.
- [50] Rajarshi Das, Shehzaad Dhuliawala, Manzil Zaheer, Luke Vilnis, Ishan Durugkar, Akshay Krishnamurthy, Alex Smola, and Andrew McCallum. Go for a walk and arrive at the answer: Reasoning over paths in knowledge bases using reinforcement learning. In *ICLR*. OpenReview.net, 2018. URL <https://openreview.net/forum?id=Syg-YfWCW>.
- [51] Jinhao Jiang, Kun Zhou, Zican Dong, Keming Ye, Xin Zhao, and Ji-Rong Wen. Structgpt: A general framework for large language model to reason over structured data. In *EMNLP*, pages 9237–9251. Association for Computational Linguistics, 2023. URL <https://doi.org/10.18653/v1/2023.emnlp-main.574>.
- [52] Jiaxin Shi, Shulin Cao, Lei Hou, Juanzi Li, and Hanwang Zhang. TransferNet: An effective and transparent framework for multi-hop question answering over relation graph. In *EMNLP*, pages 4149–4158. Association for Computational Linguistics, 2021. URL <https://aclanthology.org/2021.emnlp-main.341/>.
- [53] Xiaoxin He, Yijun Tian, Yifei Sun, Nitesh V Chawla, Thomas Laurent, Yann LeCun, Xavier Bresson, and Bryan Hooi. G-retriever: Retrieval-augmented generation for textual graph understanding and question answering. In *NeurIPS*, 2024. URL <https://openreview.net/forum?id=MPJ3oXtTZ1>.
- [54] Liyi Chen, Panrong Tong, Zhongming Jin, Ying Sun, Jieping Ye, and Hui Xiong. Plan-on-graph: Self-correcting adaptive planning of large language model on knowledge graphs. In *NeurIPS*, 2024. URL http://papers.nips.cc/paper_files/paper/2024/hash/4254e856d01a5e7b7ea050477c3ef9b9-Abstract-Conference.html.

- [55] Rajarshi Das, Ameya Godbole, Ankita Naik, Elliot Tower, Manzil Zaheer, Hannaneh Hajishirzi, Robin Jia, and Andrew McCallum. Knowledge base question answering by case-based reasoning over subgraphs. In *ICML*, volume 162 of *Proceedings of Machine Learning Research*, pages 4777–4793. PMLR, 2022. URL <https://proceedings.mlr.press/v162/das22a.html>.
- [56] Jinheon Baek, Alham Fikri Aji, and Amir Saffari. Knowledge-augmented language model prompting for zero-shot knowledge graph question answering. In *the First Workshop on Matching From Unstructured and Structured Data (MATCHING 2023)*, pages 70–98. Association for Computational Linguistics, 2023. URL <https://aclanthology.org/2023.matching-1.7/>.
- [57] Qi Zhao, Hongyu Yang, Qi Song, Xin-Wei Yao, and Xiangyang Li. Knowpath: Knowledge-enhanced reasoning via llm-generated inference paths over knowledge graphs. *CoRR*, abs/2502.12029, 2025. URL <https://doi.org/10.48550/arXiv.2502.12029>.
- [58] Yuzheng Cai, Zhenyue Guo, Yiwen Pei, Wanrui Bian, and Weiguo Zheng. Simgrag: Leveraging similar subgraphs for knowledge graphs driven retrieval-augmented generation. *CoRR*, abs/2412.15272, 2024. URL <https://doi.org/10.48550/arXiv.2412.15272>.
- [59] Jiashuo Sun, Chengjin Xu, Lumingyuan Tang, Saizhuo Wang, Chen Lin, Yeyun Gong, Lionel M. Ni, Heung-Yeung Shum, and Jian Guo. Think-on-graph: Deep and responsible reasoning of large language model on knowledge graph. In *ICLR*. OpenReview.net, 2024. URL <https://openreview.net/forum?id=nnV01PvbTv>.
- [60] Chen Zhu, Ankit Singh Rawat, Manzil Zaheer, Srinadh Bhojanapalli, Daliang Li, Felix X. Yu, and Sanjiv Kumar. Modifying memories in transformer models. *CoRR*, abs/2012.00363, 2020. URL <https://arxiv.org/abs/2012.00363>.
- [61] Yuhui Xu, Lingxi Xie, Xiaotao Gu, Xin Chen, Heng Chang, Hengheng Zhang, Zhengsu Chen, Xiaopeng Zhang, and Qi Tian. Qa-lora: Quantization-aware low-rank adaptation of large language models. In *ICLR*. OpenReview.net, 2024. URL <https://openreview.net/forum?id=WvFoJccpo8>.
- [62] Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in GPT. In *NeurIPS*, 2022. URL http://papers.nips.cc/paper_files/paper/2022/hash/6f1d43d5a82a37e89b0665b33bf3a182-Abstract-Conference.html.
- [63] Akshat Gupta, Sidharth Baskaran, and Gopala Anumanchipalli. Rebuilding ROME : Resolving model collapse during sequential model editing. In *EMNLP*, pages 21738–21744. Association for Computational Linguistics, 2024. URL <https://aclanthology.org/2024.emnlp-main.1210>.
- [64] Kevin Meng, Arnab Sen Sharma, Alex J. Andonian, Yonatan Belinkov, and David Bau. Mass-editing memory in a transformer. In *ICLR*, 2023. URL <https://openreview.net/pdf?id=MkbcAHlYgyS>.
- [65] Junfeng Fang, Houcheng Jiang, Kun Wang, Yunshan Ma, Xiang Wang, Xiangnan He, and Tat-Seng Chua. Alphaedit: Null-space constrained knowledge editing for language models. In *ICLR*, 2025. URL <https://openreview.net/forum?id=HvSyvtvg3Jh>.
- [66] Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and Baobao Chang. Can we edit factual knowledge by in-context learning? In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *EMNLP*, pages 4862–4876. Association for Computational Linguistics, 2023. URL <https://doi.org/10.18653/v1/2023.emnlp-main.296>.
- [67] Tom Hartvigsen, Swami Sankaranarayanan, Hamid Palangi, Yoon Kim, and Marzyeh Ghassemi. Aging with GRACE: lifelong model editing with discrete key-value adaptors. In *NeurIPS*, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/95b6e2ff961580e03c0a662a63a71812-Abstract-Conference.html.

- [68] Luong Trung, Xinbo Zhang, Zhanming Jie, Peng Sun, Xiaoran Jin, and Hang Li. ReFT: Reasoning with reinforced fine-tuning. In *ACL (Volume 1: Long Papers)*, pages 7601–7614. Association for Computational Linguistics, 2024. URL <https://aclanthology.org/2024.acl-long.410/>.
- [69] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering. In *EMNLP*, pages 6769–6781. Association for Computational Linguistics, 2020. URL <https://aclanthology.org/2020.emnlp-main.550/>.
- [70] Xipeng Qiu, Tianxiang Sun, Yige Xu, Yunfan Shao, Ning Dai, and Xuanjing Huang. Pre-trained models for natural language processing: A survey. *CoRR*, abs/2003.08271, 2020. URL <https://arxiv.org/abs/2003.08271>.
- [71] Xiao Ding, Bowen Chen, Li Du, Bing Qin, and Ting Liu. Cogbert: Cognition-guided pre-trained language models. In *COLING*, pages 3210–3225. International Committee on Computational Linguistics, 2022. URL <https://aclanthology.org/2022.coling-1.284>.

A Limitation and Future Work

Considering the feature space alignment, our current framework effectively reuses the original attention weights ($\mathbf{W}^Q/\mathbf{W}^K/\mathbf{W}^V$) for triple feature transformation (as elaborated in Section 3.3). However, it does not explicitly model structural information inherent in knowledge graphs. We are actively exploring integration with knowledge graph embedding (KGE) techniques through learnable transformation matrices that encode structural patterns (*e.g.*, relation hierarchies and graph connectivity). Preliminary investigations suggest that augmenting the parameter-sharing mechanism with lightweight structural adapters could enhance structure-aware representation learning while maintaining parameter efficiency. This direction may bridge the gap between discrete graph semantics and continuous language model representations, potentially improving performance on compositional reasoning tasks.