Reason-Align-Respond: Aligning LLM Reasoning with Knowledge Graphs for KGQA

Xiangqing Shen, Fanfan Wang, and Rui Xia*

School of Computer Science and Engineering, Nanjing University of Science and Technology, China {xiangqing.shen, ffwang, rxia}@njust.edu.cn

Abstract

LLMs have demonstrated remarkable capabilities in complex reasoning tasks, yet they often suffer from hallucinations and lack reliable factual grounding. Meanwhile, knowledge graphs (KGs) provide structured factual knowledge but lack the flexible reasoning abilities of LLMs. In this paper, we present Reason-Align-Respond (RAR), a novel framework that systematically integrates LLM reasoning with knowledge graphs for KGQA. Our approach consists of three key components: a Reasoner that generates human-like reasoning chains, an Aligner that maps these chains to valid KG paths, and a Responser that synthesizes the final answer. We formulate this process as a probabilistic model and optimize it using the Expectation-Maximization algorithm, which iteratively refines the reasoning chains and knowledge paths. Extensive experiments on multiple benchmarks demonstrate the effectiveness of RAR, achieving state-of-the-art performance with Hit scores of 93.3% and 91.0% on WebQSP and CWQ respectively. Human evaluation confirms that RAR generates high-quality, interpretable reasoning chains well-aligned with KG paths. Furthermore, RAR exhibits strong zero-shot generalization capabilities and maintains computational efficiency during inference.

1 Introduction

Large language models (LLMs) have exhibited impressive capabilities across a range of complex tasks (Hadi et al., 2023), yet their reasoning processes often lack reliable factual knowledge. This shortcoming reduces interpretability, leads to hallucinations, and causes factual or logical errors (Huang et al., 2025). Knowledge graphs (KGs) (Bollacker et al., 2008), which organize factual knowledge in a structured format, offer strong interpretability and expressive power, making them

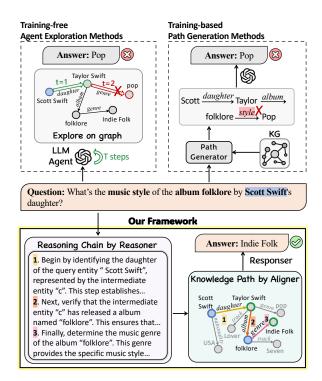


Figure 1: The comparison between our Reason-Align-Respond framework and the existing methods for LLM-based KGQA.

a reliable source of factual support for LLM reasoning. Integrating KGs with LLMs in question answering, *e.g*, knowledge graph question answering (KGQA), has gained interest as an effective strategy to mitigate hallucinations and enhance interpretability.

The existing LLM-based KGQA studies broadly fall into two main categories: Training-free Agent Exploration methods (Sun et al., 2024) and Training-based Path Generation methods (Luo et al., 2024a). The former uses LLMs as agents to explore nodes in KGs, while the latter trains generators to retrieve and generate knowledge paths. However, both methods lack the holistic reasoning process that humans typically employ when answering questions. This limitation, combined

^{*}Corresponding author.

with the semantic gap between natural language descriptions and structured knowledge graphs, often results in reasoning paths that do not align with human logic, lack semantic relevance, or contain unnecessary noise, as shown in Fig. 1.

On the other hand, the latest deep reasoning LLMs, such as OpenAI-o1 (Zhong et al., 2024) and DeepSeek-R1 (Guo et al., 2025), produce holistic reasoning chains before answering challenging questions, showcasing stronger logical reasoning abilities. But the reasoning chains obtained through reinforcement learning tend to be verbose, sometimes contain logical fallacies, and amplify hallucinations (Bao et al., 2025). While our work is not specifically aimed at improving Deepseek-R1, we believe that using knowledge paths from KGs as constraints for deep reasoning might potentially alleviate these issues.

The two aforementioned aspects are complementary. For one thing, allowing LLMs to first perform human-like reasoning in KGQA, can establish a structured, goal-oriented framework that helps reduce invalid or noisy KG path exploration. For another, while free reasoning by LLMs tends to increase hallucinations, incorporating knowledge paths from KGs provides constraints that help ensure alignment with the knowledge graph, thereby reducing hallucinations.

In this work, we introduce Reason-Align-Respond (RAR), a novel framework that systematically integrates three modules—Reasoner, Aligner, and Responser—to align LLM reasoning with knowledge graphs for KGQA. Each of the three modules is a fine-tuned LLM. Firstly, Reasoner conducts global reasoning for the question, simulating human-like thinking process to generate a reasoning chain that guides LLM's exploration on the KG. Secondly, Aligner decodes a knowledge path on the KG based on the above reasoning chain. Each decoding step ensures correspondence to an actual knowledge triple in the graph. Finally, Responser leverages both the reasoning chain and the knowledge path to generate the final answer.

We treat both reasoning chain z_r and knowledge path z_p as latent variables, use a probabilistic model to formalize the distribution of answer a given the question q and KG \mathcal{G} , and propose an end-to-end training algorithm to jointly fine-tune the parameters across all three modules. Specifically, we employ the expectation-maximization (EM) algorithm to iteratively optimize model parameters while inferring the latent variables. The

E-step estimates the posterior probability of latent variables (z_r, z_p) given observations (q, \mathcal{G}, a) and samples high-quality reasoning chain and aligned knowledge paths accordingly. The M-step maximizes the evidence lower bound (ELBO) of the log-likelihood to update model parameters. Through iterative optimization, the model progressively refines the reasoning chain and knowledge path, and in turn promotes the generation of high-quality responses, ultimately forming a stable closed loop.

We conduct extensive evaluation across multiple benchmarks, including WebQSP, CWQ, CommonsenseQA (CSQA), and MedQA, using Freebase, ConceptNet and a medical KG as knowledge graphs. The results demonstrate the effectiveness of RAR in improving KGQA performance. Compared to 19 baselines across three categories, RAR achieves state-of-the-art results. Specifically, it achieves Hit scores of 93.3% on WebQSP and 91.0% on CWQ, with corresponding F1 score of 87.7% and 84.8%, significantly surpassing the existing methods. Additionally, RAR demonstrates strong zero-shot generalizability on CSQA and MedQA. The EM algorithm demonstrates gradual performance improvements with increasing iterations and shows stable convergence after 200 steps. Human evaluation of reasoning chains shows that RAR generates high-quality, interpretable chains aligned with KG paths, ensuring accurate and reliable reasoning. We also report the results of RAR under different LLM backbone models, and conduct the ablation study that confirms the importance of each component. Notably, RAR achieves these effective performance while maintaining computational efficiency during inference. These comprehensive results demonstrate that our approach successfully bridges the gap between LLM reasoning and structured knowledge, offering a promising direction for developing more reliable and interpretable question answering systems.

2 Approach

Knowledge graphs (KGs) contain extensive factual knowledge in the form of triples: $\mathcal{G} = \{(e,r,e')|e,e'\in\mathcal{E},r\in\mathcal{R}\}$, where \mathcal{E} and \mathcal{R} denote sets of entities and relations, respectively. The goal of knowledge graph question answering (KGQA) is to retrieve relevant facts from a KG \mathcal{G} and generate an answer a in response to a given natural language question q.

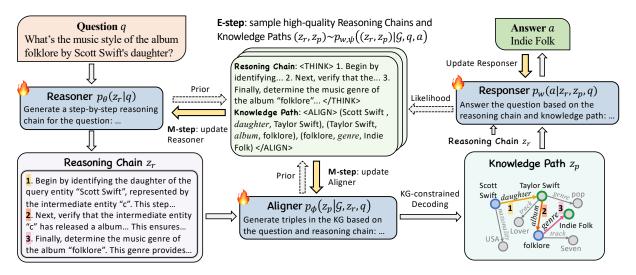


Figure 2: Illustration of our RAR framework comprising Reasoner, Aligner, Responser with iterative EM optimization.

2.1 Task Formalization

We introduce two latent variables, a reasoning chain z_r and a knowledge path z_p , working together to answer the question q based on a KG \mathcal{G} :

- Reasoning Chain z_r denotes a chain of discrete reasoning steps expressed in natural language, working together to address the question q.
- Knowledge Path z_p denotes an interconnected path of knowledge triples extracted from a KG G.

Since neither z_r nor z_p are explicitly annotated in existing KGQA benchmarks, we treat them as latent variables, and employ a probabilistic model to formalize the distribution of the answer a conditioned on the question q and KG \mathcal{G} as:

$$p_{w,\phi,\theta}(a|\mathcal{G},q) = \sum_{z_r,z_p} p_w(a|z_r,z_p,q) p_\phi(z_p|\mathcal{G},z_r,q) p_\theta(z_r|q),$$
(1)

where we assume a is conditionally independent of \mathcal{G} given (z_r, z_p, q) , allowing factorization and summation over these conditional probabilities. On this basis, we introduce our framework, Reason-Align-Respond (RAR), that integrates three modules—Reasoner, Aligner, and Responser—to align LLM reasoning with knowledge graphs for KGQA, as illustrated in Fig. 2.

It is worth noting that each of the three modules is a fine-tuned LLM, with parameters denoted as θ , ϕ , and ω , respectively. They each utilize the Prompts shown in App. F to generate reasoning chains, knowledge paths, and answers.

Firstly, **Reasoner** $p_{\theta}(z_r|q)$ generates a latent reasoning chain z_r to address the question q. The Reasoner generates a reasoning chain in the following format:

$$z_r = \langle THINK \rangle s_1 \dots s_t \langle THINK \rangle$$
,

where <THINK> and </THINK> are special tokens denoting the start and end of the reasoning process, and each s_i is a discrete reasoning step.

Secondly, **Aligner** $p_{\phi}(z_p|\mathcal{G}, z_r, q)$ takes the question q and the reasoning chain z_r as inputs to explore KG \mathcal{G} , producing a latent reasoning path z_p that aligns with z_r . Aligner processes the prompt with q and z_r to generate a knowledge path in the following format:

$$z_p = \langle ALIGN \rangle \pi_1 \dots \pi_t \langle ALIGN \rangle$$

where <ALIGN> and </ALIGN> mark the beginning and end of the knowledge path. Each π_i is a triple from the KG $\mathcal G$ formatted as $\pi_i = \langle \mathsf{TRI} \rangle (e_i^h, r_i, e_i^t) \langle \mathsf{TRI} \rangle$, where $\langle \mathsf{TRI} \rangle$ and $\langle \mathsf{TRI} \rangle$ bound the triple.

Finally, **Responser** $p_w(a|z_r, z_p, q)$ generates the final answer a by synthesizing the question q, reasoning chain z_r , and knowledge path z_p .

2.2 Optimization via the EM Algorithm

As we have mentioned, each of the three modules is a fine-tuned LLM. In this section, we propose an end-to-end training algorithm to jointly optimize the parameters across all three modules. The training objective is the likelihood of the distribution $p_{w,\phi,\theta}(a|\mathcal{G},q)$ shown in Eq. (1). Since our model involves latent variables z_r and z_p , we adopt

the Expectation-Maximization (EM) algorithm—a principled approach for maximum likelihood estimation (MLE) in latent-variable models (Dempster et al., 1977; Sen et al., 2022; Qu et al., 2021).

Unifying Reasoner and Aligner. In practice, we unify Reasoner and Aligner by merging their latent variables z_r and z_p into a single one $z=(z_r,z_p)$. This results in a consolidated module, referred to as **ReAligner**, whose parameters denoted by ψ . Hence, instead of generating z_r and z_p separately, ReAligner simultaneously outputs both a Reasoning Chain and a Knowledge Path, treating it as a single instruction-tuning task with a prompt template (see App. F). With this simplification, the conditional distribution for the final answer a is:

$$p_{w,\psi}(a|\mathcal{G},q) = \sum_{z} p_w(a|q,z) p_{\psi}(z|\mathcal{G},q), \quad (2)$$

where \mathcal{G} denotes the KG, q the question, and z the unified latent variable (combining the Reasoning Chain and Knowledge Path).

Learning Objective. We aim to learn the parameters (w, ψ) by maximizing the log-likelihood of the training data with respect to Eq. (2), written as:

$$\max_{w,\psi} \mathcal{O}(w,\psi) = \log \mathbb{E}_{z \sim p_{\psi}(z|\mathcal{G},q)}[p_w(a|q,z)].$$
(3)

According to Jensen's inequality, we have:

$$\mathcal{O}(w,\psi) \ge \underbrace{\mathbb{E}_{q(z)} \log(\frac{p_w(a|q,z)p_{\psi}(z|\mathcal{G},q)}{q(z)})}_{\mathcal{L}_{\text{ELBO}}},$$

where q(z) is a variational distribution. Equality holds when $q(z) = p_{w,\psi}(z|\mathcal{G},q,a)$, the true posterior of z. The term $\mathcal{L}_{\text{ELBO}}$ is the Evidence Lower BOund (ELBO), and maximizing $\mathcal{L}_{\text{ELBO}}$ indirectly maximizes $\mathcal{O}(w,\psi)$.

EM Algorithm. The EM algorithm alternates between an E-step and an M-step until convergence:

E-step. Given current parameters $(w^{(t)}, \psi^{(t)})$ at iteration t, E-step updates the variational distribution $q^{(t)}(z)$ by minimizing $\mathrm{KL}(q(z)||p_{w^{(t)},\psi^{(t)}}(z|\mathcal{G},q,a))$. The solution is the posterior of z under the current parameters:

$$q^{(t)}(z) = p_{w(t) | y/(t)}(z|\mathcal{G}, q, a).$$
 (5)

M-step. Keeping $q^{(t)}(z)$ fixed, M-step maximizes $\mathcal{L}_{\text{ELBO}}$ in Eq. (4) with respect to w and ψ .

Ignoring terms that do not depend on (w, ψ) , the objective reduces to:

$$\begin{split} Q(w, \psi | w^{(t)}, \psi^{(t)}) &= \sum_{(\mathcal{G}, q, a)} \sum_{z} q^{(t)}(z) \, \log[p_w(a|q, z) p_\psi(z|\mathcal{G}, q)] \\ &= \underbrace{\sum_{(\mathcal{G}, q, a)} \sum_{z} q^{(t)}(z) \log p_w(a|q, z)}_{Q_{\text{Responser}}(w)} \\ &+ \underbrace{\sum_{(\mathcal{G}, q, a)} \sum_{z} q^{(t)}(z) \log p_\psi(z|\mathcal{G}, q)}_{Q_{\text{ReAligner}}(\psi)} \end{split}$$

which naturally divides into the instruction-tuning objective for Responser and ReAligner in Eq. (2).

By iterative optimization with the EM algorithm, our framework progressively refines its understanding of the question. This iterative process gradually corrects any flaws in Reasoning Chains and Knowledge Paths, leading to answers that are both higher in quality and more interpretable, while significantly reducing the risk of hallucination.

The workflow of the EM algorithm is shown in Alg. 1, with more details in practice in App. A.

Algorithm 1 The EM algorithm in RAR

while not converge do

For each instance, sample N Reasoning Chains and Knowledge Paths z_I from Re-Aligner p_{ψ} .

For each instance, update Responser p_w with $Q_{\text{Responser}}(w)$ in Eq. (6) using z_I .

 \Box *E-step:*

For each instance, identify K high-quality Reasoning Chains and Knowledge Paths z_I^h from z_I based on Eq. (5).

 \square *M-step*:

For each instance, update ReAligner p_{ψ} according to $Q_{\text{ReAligner}}(\psi)$ in Eq. (6).

end while

2.3 Techniques During Inference

During inference, given q, Reasoner generates z_r , while Aligner produces z_p . Responser synthesizes them to produce a. To enhance performance, we introduce three additional key techniques.

KG-constrained Decoding. KG-constrained Decoding aims to prevent hallucinated triples that do

Types	Methods		WebQSP		CWQ	
1) pes			F1	Hit	F1	
	Qwen2-7B (Yang et al., 2024)	50.8	35.5	25.3	21.6	
	Llama-2-7B (Touvron et al., 2023)	56.4	36.5	28.4	21.4	
	Llama-3.1-8B (Meta, 2024)	55.5	34.8	28.1	22.4	
	GPT-4o-mini (OpenAI, 2024a)	63.8	40.5	63.8	40.5	
LLM Reasoning	ChatGPT (OpenAI, 2022)	59.3	43.5	34.7	30.2	
LLIVI Reasoning	ChatGPT+Few-shot (Brown et al., 2020)	68.5	38.1	38.5	28.0	
	ChatGPT+CoT (Wei et al., 2022b)	73.5	38.5	47.5	31.0	
	ChatGPT+Self-Consistency (Wang et al., 2024)	83.5	63.4	56.0	48.1	
	GraftNet (Sun et al., 2018)	66.7	62.4	36.8	32.7	
Cronh Daggaring	NSM (He et al., 2021)	68.7	62.8	47.6	42.4	
Graph Reasoning	SR+NSM (Zhang et al., 2022)	68.9	64.1	50.2	47.1	
	ReaRev (Mavromatis and Karypis, 2022)	76.4	70.9	52.9	47.8	
	KD-CoT (Wang et al., 2023a)	68.6	52.5	55.7	-	
	EWEK-QA (Dehghan et al., 2024)	71.3	-	52.5	-	
	ToG (GPT-4) (Sun et al., 2024)	82.6	-	68.5	-	
KG+LLM	EffiQA (Dong et al., 2025)	82.9	-	69.5		
	RoG (Llama-2-7B) (Luo et al., 2024b)	85.7	70.8	62.6	56.2	
	GNN-RAG+RA (Mavromatis and Karypis, 2024)	90.7	73.5	68.7	60.4	
	GCR (Llama-3.1-8B + GPT-4o-mini) (Luo et al., 2024c)	92.2	74.1	75.8	61.7	
	RAR (Llama-3.1-8B + GPT-4o-mini)	93.3	87.7	91.0	84.8	

Table 1: Performance comparison with different baselines on the two KGQA datasets.

not exist in the KG. When generating the Knowledge Path, Aligner may inadvertently produce triples absent from the KG. To address this, KG-constrained Decoding restricts the output tokens so that only tokens forming valid KG triples can be produced. In this way, the generated Knowledge Path strictly aligns with actual entities and relations in the KG. Related work (Luo et al., 2024c; Li et al., 2024) also attempts to mitigate similar issues; our approach is tailored specifically to our framework.

Knowledge Path Expansion. Knowledge Path Expansion addresses the potential incompleteness of initially generated Knowledge Paths. To illustrate, consider a question about countries that share borders with the United States. a Knowledge Path

<ALIGN><TRI>(US, borders, Mexico)
is generated by Aligner. While correct, this represents only one instance of a broader pattern. By abstracting the specific instance into a template:
<ALIGN><TRI>(US, borders, ?country)</TRI></ALIGN>, where ?country is a variable, we capture the fundamental relationship structure. Applied to the KG,

this template retrieves all valid instances, such as:

<ALIGN><TRI>(US, borders, Canada)
/TRI></ALIGN>.
This method transforms a single Knowledge Path into a comprehensive query template, enabling more complete and exhaustive answers.

LLM-driven Consolidation. LLM-driven Consolidation addresses the challenge of inconsistencies and noise that emerge when sampling multiple Reasoning Chains and Knowledge Paths. Multiple sampling helps increase coverage and improve the likelihood of correct answers, but inevitably introduces noise and conflicts between samples. To address this challenge, we propose using a powerful LLM as a "Consolidator" that analyzes and integrates multiple Reasoning Chains and Knowledge Paths to derive final answers, following the prompt template detailed in App. F. This approach effectively preserves the benefits of multiple sampling while leveraging the LLM's analytical capabilities to produce reliable answers.

3 Experiment

3.1 Experiment Settings

Datasets. Following previous research (Luo et al., 2024c; Sun et al., 2024), we evaluate our model on three datasets: WebQuestionSP (WebQSP)) (Yih et al., 2016), Complex WebQuestions (CWQ) (Talmor and Berant, 2018), and CommonsenseQA (CSQA) (Talmor et al., 2019). The first two datasets use Freebase (Bollacker et al., 2008) as the KG, while CSQA leverages ConceptNet (Speer et al., 2017), allowing us to assess model generalizability across the unseen KG.

Baselines. We compare RAR with 19 baselines across three categories: LLM reasoning methods, graph reasoning methods, and KG-enhanced LLM reasoning methods.

Evaluation Metrics. For WebQSP and CWQ, we adopt Hit and F1 as evaluation metrics. Hit checks whether the generated predictions match any correct answer, while F1 evaluates overall answer coverage by balancing precision and recall. For CSQA, a multiple-choice QA dataset, we use accuracy as the evaluation metric.

Implementations. Our implementation uses Llama-3.1-8B (Meta, 2024) as the backbone for Reasoner, Aligner, and Responser. To enhance question decomposition, we pretrain both Reasoner and Aligner using 2,000 exemplars demonstrating step-by-step KG-based problem-solving. For each component, we generate top-10 candidates using KG-constrained Decoding and Knowledge Path Expansion, with GPT-40-mini handling LLM-driven Consolidation. Details are provided in App. C.

3.2 Main Results

Tab. 1 shows that RAR achieves significant gains on both WebQSP and CWQ datasets, improving the state-of-the-art by 13.6% and 23.1% in F1, and by 1.1% and 15.2% in Hit respectively. These remarkable improvements demonstrate the effectiveness of integrating structured KG knowledge with LLM reasoning.

3.3 Ablation Study

As shown in Tab. 2, removing LLM-driven Consolidation (LC) lowers precision but increases recall, since LC aims to eliminate noisy predictions. Excluding Reasoner causes a pronounced drop in precision but a rise in recall, indicating that Reasoning Chains guide Aligner to explore the KG more

F1	Precision	Recall
84.8	84.9	89.0
83.0	81.8	91.2
80.3	75.6	94.7
48.9	54.3	50.8
71.7	80.2	71.7
	84.8 83.0 80.3 48.9	84.8 84.9 83.0 81.8 80.3 75.6 48.9 54.3

Table 2: Impact of different components on CWQ.

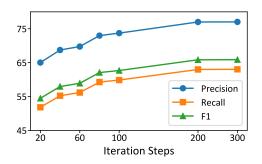


Figure 3: Impact of iteration steps of the EM algorithm.

accurately, reducing hallucination and noise. Disabling Knowledge Path Expansion (KPE) diminishes performance, confirming its role in enriching Knowledge Paths. Most importantly, removing KG-constrained Decoding (KD) yields the largest performance decrease, underscoring the importance of restricting generation to valid KG paths.

3.4 Further Analyses

Impact of Iteration Steps. As shown in Fig. 3, RAR exhibits consistent improvement across all metrics during EM updates. The performance rises rapidly in early stages and continues to refine over iterations, eventually reaching convergence with minimal fluctuations.

Quality of Reasoning Chains. Through manual evaluation of 500 randomly selected samples, we assess the quality of Reasoning Chains using two criteria: reasoning correctness and KG alignment. The correctness metric evaluates whether the Reasoning Chain successfully solves the given question, while the alignment metric measures how well the reasoning steps correspond to valid KG paths. As shown in Fig. 4, RAR substantially outperforms both GPT-40 and baseline methods across both metrics. These results demonstrate that by aligning Reasoning Chains to KG structures, our approach not only improves the reliability of the reasoning process but also enhances its interpretability.

KG-constrained Decoding Effectiveness. We examine the effectiveness of KG-constrained Decod-

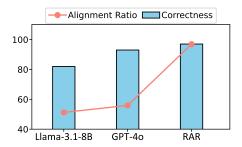


Figure 4: Human evaluation of reasoning chains on CWQ.

Types	Methods	Hit	Avg. Runtime (s)
Path Generation	GNN-RAG RoG GCR	66.8 62.6 75.8	1.73 2.68 3.72
Agent Exploration	ToG EffiQA	68.5 69.5	18.89
Ours	RAR	91.0	4.38

Table 3: Efficiency and performance of RAR compared to different methods on CWQ.

ing in mitigating hallucinations and maintaining efficiency. Our method achieves zero hallucinations in Knowledge Paths when answers are correct, while without constraints, even correct answers show a 44% hallucination rate. The efficiency evaluation reveals minimal computational overhead. Particularly noteworthy is the comparison with GCR, the previous state-of-the-art method using KG constraints. Our approach achieves a 15.2% improvement in answer accuracy over GCR, with only a marginal increase in runtime from Reasoning Chain generation. This modest overhead is well justified by the significant gains in answer reliability and interpretability.

Impact of Beam Size. As shown in Fig. 5, increasing beam size allows RAR to explore more potential Reasoning Chains and Knowledge Paths, leading to improved performance across all metrics. This demonstrates that examining multiple candidate solutions helps identify better Reasoning Chains and Knowledge Paths for responses of higher quality.

Impact of different LLM Backbone. Tab. 4 demonstrates that larger LLMs generally achieve better performance, with Llama-3.1-8B and GPT-40 delivering the strongest results for backbone and LLM-based Consolidation (LC), respectively.

Zero-shot Generalizability to Unseen KGs. Following Luo et al. (2024c), we evaluate RAR's zero-shot transfer capabilities on CSQA (Talmor et al.,

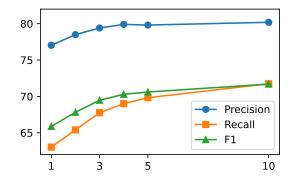


Figure 5: Impact of different beam size on CWQ.

Components	Variants	Hit	F1
	Llama-2-7B	84.0	68.9
Reasoner and	Llama-3.1-8B	85.2	72.6
Troubonion unio	Qwen2-0.5B	71.3	56.0
Aligner	Qwen2-1.5B	72.0	56.6
	Qwen2-7B	81.4	67.1
	GPT-4o-mini	91.0	84.8
LLM-driven	GPT-40	92.8	84.9
Consolidation	Qwen2-7B	88.7	82.2
Consolidation	Llama-3.1-8B	90.6	83.7
	Llama-3.1-70B	92.4	83.0

Table 4: Impact of using different LLM backbones for Reasoner, Aligner and LLM-driven Consolidation.

2019) and MedQA (Jin et al., 2021). The results show that RAR achieves superior zero-shot performance compared to GPT-40-mini on both datasets. Notably, RAR achieves comparable performance to GCR, as both methods leverage the KG to constrain the decoding process to enhance reasoning without requiring additional training on the target KG.

3.5 Case study

Fig. 6 illustrates the qualitative results by our RAR framework. To investigate the effect of EM iterations, in Fig. 7, we also present two representative cases that demonstrate the evolution of RAR's behavior across iterations of the EM algorithm. These cases showcase distinct improvements in RAR's ability to generate Reasoning Chains and Knowledge Paths. In Case 1, we examine RAR's response to the query "What high school did the artist who recorded 'Girl Tonight' attend?" The early-stage model demonstrates suboptimal verification behavior by directly searching for high school information without following the proper verification sequence: first identifying educational institutions, then verifying the specific type as high school. This Reasoning Chain prevents proper alignment with

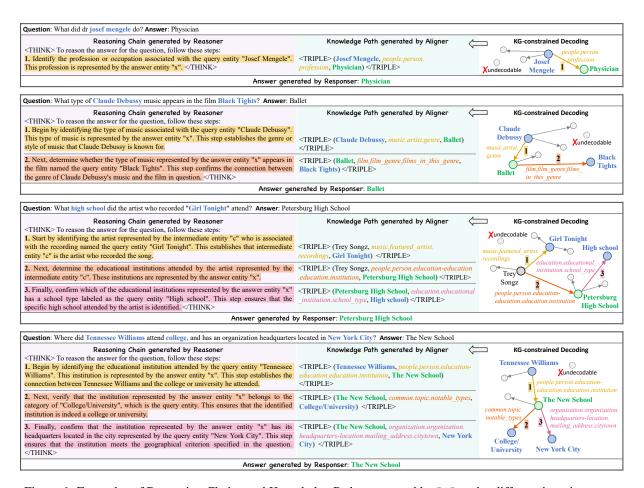


Figure 6: Examples of Reasoning Chains and Knowledge Paths generated by RAR under different iteration steps.

Model	CSQA	MedQA	
GPT-4o-mini	91	75	
GCR	94	79	
RAR	94	80	

Table 5: Zero-shot transferability to unseen KG.

the Knowledge Path in the KG. In contrast, the laterstage model generates both Reasoning Chains and Knowledge Paths effectively, producing a higher quality Reasoning Chain that successfully aligns with the Knowledge Path in the KG. Case 2 examines RAR's handling of "What type of Claude Debussy music appears in the film Black Tights?" Here, we observe a different pattern of improvement. While the early-stage model generates the same Reasoning Chain as the later-stage model, it fails to generate Knowledge Paths that fully align with and reflect this reasoning, resulting in a misaligned Knowledge Path that do not lead to the correct answer. The later-stage model maintains consistency between Reasoning Chains and Knowledge Paths, thus arriving at the correct answer. These cases validate the effectiveness of the EM approach.

4 Related Work

LLM Reasoning. Recent progress in LLMs have spurred extensive research to improve deep reasoning. One line of work focuses on promptingbased methods, which elicit intermediate reasoning steps during inference, such as Chain-of-Thought (CoT) prompting (Wei et al., 2022a). Building on CoT, self-consistency (Wang et al., 2023b) generates multiple reasoning traces and selects the most consistent answer. Tree-of-Thought (Yao et al., 2023) explores branching steps in a tree structure to uncover optimal solutions. Beyond prompting, researchers have investigated fine-tuning LLMs on reasoning tasks (Yu et al., 2022; Hoffman et al., 2023), including reinforcement learning methods (OpenAI, 2024b; Guo et al., 2025) that encourage more complex multi-step reasoning before arriving at a final answer.

KG-enhanced LLM Reasoning. Despite their remarkable performance, LLMs still face limitations

such as incomplete or outdated domain knowledge, interpretability challenges, and the potential for hallucinations. To address these issues, a growing body of work has focused on integrating LLMs with KGs (Pan et al., 2024). KD-CoT (Wang et al., 2023a) enhances CoT by retrieving relevant facts from external KGs, guiding LLMs with more reliable information. RoG (Luo et al., 2024b) employs a plan-retrieve-reason pipeline that explicitly fetches KG-based reasoning paths to ground the final answer. GCR (Luo et al., 2024c) further mitigates hallucinations by enforcing graphconstrained decoding, ensuring that every reasoning step aligns with valid KG connections. GNN-RAG (Mavromatis and Karypis, 2024) leverages a graph neural network for effective KG retrieval, while StructGPT (Jiang et al., 2023) and ToG (Sun et al., 2024) treat the LLM as an agent that navigates the KG, assembling multi-hop paths to produce more trustworthy answers.

5 Conclusion

In this paper, we present RAR, a novel framework that integrates LLM reasoning with knowledge graphs for KGQA through three key components—Reasoner, Aligner, and Responser. We formulate this process as a probabilistic model and optimize it using the Expectation-Maximization algorithm. Through extensive experiments on multiple benchmarks, we demonstrate that RAR achieves state-of-the-art performance while maintaining interpretability and computational efficiency. These results demonstrate that RAR successfully bridges the gap between LLM reasoning and structured knowledge, offering a promising direction for building reliable and interpretable QA systems.

Limitations

One limitation of our framework lies in the computational overhead introduced by Reasoner. In certain cases, especially for complex queries, the Reasoning Chain generated by Reasoner can become relatively long, increasing resource consumption. However, our experimental results demonstrate that the performance gains from incorporating Reasoning Chain justify this additional cost, striking a practical balance between efficiency and effectiveness. Another limitation concerns the generalizability to specialized domains. Though our framework, trained on Freebase-based KGQA datasets, shows improved generalization to unseen KGs compared

to previous methods, its performance on highly specialized knowledge graphs (*e.g.*, medical KGs) remains to be enhanced. Improving adaptation to domain-specific KGs presents a promising direction for future research.

References

Forrest Bao, Chenyu Xu, and Ofer Mendelevitch. 2025. Deepseek-r1 hallucinates more than deepseek-v3. Accessed: 2025-02-10.

Kurt D. Bollacker, Colin Evans, Praveen K. Paritosh, Tim Sturge, and Jamie Taylor. 2008. Freebase: a collaboratively created graph database for structuring human knowledge. In *Proceedings of the ACM SIG-MOD International Conference on Management of Data, SIGMOD 2008, Vancouver, BC, Canada, June 10-12*, 2008, pages 1247–1250. ACM.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901.

Sitao Cheng, Ziyuan Zhuang, Yong Xu, Fangkai Yang, Chaoyun Zhang, Xiaoting Qin, Xiang Huang, Ling Chen, Qingwei Lin, Dongmei Zhang, Saravan Rajmohan, and Qi Zhang. 2024. Call me when necessary: LLMs can efficiently and faithfully reason over structured environments. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 4275–4295, Bangkok, Thailand. Association for Computational Linguistics.

Mohammad Dehghan, Mohammad Alomrani, Sunyam Bagga, David Alfonso-Hermelo, Khalil Bibi, Abbas Ghaddar, Yingxue Zhang, Xiaoguang Li, Jianye Hao, Qun Liu, Jimmy Lin, Boxing Chen, Prasanna Parthasarathi, Mahdi Biparva, and Mehdi Rezagholizadeh. 2024. EWEK-QA: Enhanced web and efficient knowledge graph retrieval for citation-based question answering systems. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14169–14187, Bangkok, Thailand. Association for Computational Linguistics.

Arthur P Dempster, Nan M Laird, and Donald B Rubin. 1977. Maximum likelihood from incomplete data via the em algorithm. *Journal of the royal statistical society: series B (methodological)*, 39(1):1–22.

Zixuan Dong, Baoyun Peng, Yufei Wang, Jia Fu, Xiaodong Wang, Xin Zhou, Yongxue Shan, Kangchen Zhu, and Weiguo Chen. 2025. Effiqa: Efficient question-answering with strategic multi-model collaboration on knowledge graphs. In *Proceedings of the 31st International Conference on Computational Linguistics, COLING 2025, Abu Dhabi, UAE, January 19-24, 2025*, pages 7180–7194. Association for Computational Linguistics.

- Yu Gu, Yiheng Shu, Hao Yu, Xiao Liu, Yuxiao Dong, Jie Tang, Jayanth Srinivasa, Hugo Latapie, and Yu Su. 2024. Middleware for LLMs: Tools are instrumental for language agents in complex environments. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 7646–7663, Miami, Florida, USA. Association for Computational Linguistics.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Muhammad Usman Hadi, Rizwan Qureshi, Abbas Shah, Muhammad Irfan, Anas Zafar, Muhammad Bilal Shaikh, Naveed Akhtar, Jia Wu, Seyedali Mirjalili, et al. 2023. A survey on large language models: Applications, challenges, limitations, and practical usage. *Authorea Preprints*, 3.
- Gaole He, Yunshi Lan, Jing Jiang, Wayne Xin Zhao, and Ji-Rong Wen. 2021. Improving multi-hop knowledge base question answering by learning intermediate supervision signals. In *Proceedings of the 14th ACM international conference on web search and data mining*, pages 553–561.
- Matthew Douglas Hoffman, Du Phan, David Dohan, Sholto Douglas, Tuan Anh Le, Aaron Parisi, Pavel Sountsov, Charles Sutton, Sharad Vikram, and Rif A. Saurous. 2023. Training chain-of-thought via latent-variable inference. In Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 16, 2023.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2025. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems*, 43(2):1–55.
- Xiang Huang, Sitao Cheng, Shanshan Huang, Jiayu Shen, Yong Xu, Chaoyun Zhang, and Yuzhong Qu. 2024. QueryAgent: A reliable and efficient reasoning framework with environmental feedback based self-correction. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5014–5035, Bangkok, Thailand. Association for Computational Linguistics.
- Jinhao Jiang, Kun Zhou, Zican Dong, Keming Ye, Wayne Xin Zhao, and Ji-Rong Wen. 2023. Structgpt: A general framework for large language model to reason over structured data. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9237–9251.
- Jinhao Jiang, Kun Zhou, Xin Zhao, and Ji-Rong Wen. 2022. Unikgqa: Unified retrieval and reasoning for

- solving multi-hop question answering over knowledge graph. In *The Eleventh International Conference on Learning Representations*.
- Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2021. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied Sciences*, 11(14):6421.
- Kun Li, Tianhua Zhang, Xixin Wu, Hongyin Luo, James Glass, and Helen Meng. 2024. Decoding on graphs: Faithful and sound reasoning on knowledge graphs through generation of well-formed chains. *arXiv* preprint arXiv:2410.18415.
- Tianle Li, Xueguang Ma, Alex Zhuang, Yu Gu, Yu Su, and Wenhu Chen. 2023. Few-shot in-context learning on knowledge base question answering. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6966–6980, Toronto, Canada. Association for Computational Linguistics.
- Linhao Luo, Yuan-Fang Li, Gholamreza Haffari, and Shirui Pan. 2024a. Reasoning on graphs: Faithful and interpretable large language model reasoning. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.* OpenReview.net.
- Linhao Luo, Yuan-Fang Li, Gholamreza Haffari, and Shirui Pan. 2024b. Reasoning on graphs: Faithful and interpretable large language model reasoning. In *International Conference on Learning Representations*
- Linhao Luo, Zicheng Zhao, Chen Gong, Gholamreza Haffari, and Shirui Pan. 2024c. Graph-constrained reasoning: Faithful reasoning on knowledge graphs with large language models. *CoRR*, abs/2410.13080.
- Costas Mavromatis and George Karypis. 2022. Rearev: Adaptive reasoning for question answering over knowledge graphs. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2447–2458.
- Costas Mavromatis and George Karypis. 2024. Gnnrag: Graph neural retrieval for large language model reasoning. *arXiv preprint arXiv:2405.20139*.
- Meta. 2024. Build the future of ai with meta llama 3.
- Zhijie Nie, Richong Zhang, Zhongyuan Wang, and Xudong Liu. 2024. Code-style in-context learning for knowledge-based question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 18833–18841.
- OpenAI. 2022. Introducing chatgpt.
- OpenAI. 2024a. Hello gpt-4o.
- OpenAI. 2024b. Learning to reason with llms.

- Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Jiapu Wang, and Xindong Wu. 2024. Unifying large language models and knowledge graphs: A roadmap. *IEEE Transactions on Knowledge and Data Engineering (TKDE)*.
- Meng Qu, Junkun Chen, Louis-Pascal A. C. Xhonneux, Yoshua Bengio, and Jian Tang. 2021. Rnnlogic: Learning logic rules for reasoning on knowledge graphs. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net.
- Prithviraj Sen, Breno W. S. R. de Carvalho, Ryan Riegel, and Alexander G. Gray. 2022. Neuro-symbolic inductive logic programming with logical neural networks. In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI 2022, The Twelveth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 March 1, 2022*, pages 8212–8219. AAAI Press.
- Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. In *Proceedings of the AAAI conference on artificial intelligence*, volume 31.
- Haitian Sun, Bhuwan Dhingra, Manzil Zaheer, Kathryn Mazaitis, Ruslan Salakhutdinov, and William Cohen. 2018. Open domain question answering using early fusion of knowledge bases and text. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4231–4242.
- Jiashuo Sun, Chengjin Xu, Lumingyuan Tang, Saizhuo Wang, Chen Lin, Yeyun Gong, Lionel Ni, Heung-Yeung Shum, and Jian Guo. 2024. Think-on-graph: Deep and responsible reasoning of large language model on knowledge graph. In *The Twelfth International Conference on Learning Representations*.
- Alon Talmor and Jonathan Berant. 2018. The web as a knowledge-base for answering complex questions. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 641–651.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. Commonsenseqa: A question answering challenge targeting commonsense knowledge. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

- Keheng Wang, Feiyu Duan, Sirui Wang, Peiguang Li, Yunsen Xian, Chuantao Yin, Wenge Rong, and Zhang Xiong. 2023a. Knowledge-driven cot: Exploring faithful reasoning in llms for knowledge-intensive question answering. *arXiv preprint arXiv:2308.13259*.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023b. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023.* OpenReview.net.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2024. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022a. Chain-of-thought prompting elicits reasoning in large language models. In Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 December 9, 2022.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022b. Chain-of-thought prompting elicits reasoning in large language models. Advances in Neural Information Processing Systems, 35:24824–24837.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. 2024. Qwen2 technical report. arXiv preprint arXiv:2407.10671.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. In Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 16, 2023.
- Wen-tau Yih, Matthew Richardson, Christopher Meek, Ming-Wei Chang, and Jina Suh. 2016. The value of

semantic parse labeling for knowledge base question answering. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 2: Short Papers.* The Association for Computer Linguistics.

Ping Yu, Tianlu Wang, Olga Golovneva, Badr AlKhamissy, Gargi Ghosh, Mona T. Diab, and Asli Celikyilmaz. 2022. ALERT: adapting language models to reasoning tasks. *CoRR*, abs/2212.08286.

Jing Zhang, Xiaokang Zhang, Jifan Yu, Jian Tang, Jie Tang, Cuiping Li, and Hong Chen. 2022. Subgraph retrieval enhanced model for multi-hop knowledge base question answering. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5773–5784.

Tianyang Zhong, Zhengliang Liu, Yi Pan, Yutong Zhang, Yifan Zhou, Shizhe Liang, Zihao Wu, Yanjun Lyu, Peng Shu, Xiaowei Yu, Chao Cao, Hanqi Jiang, Hanxu Chen, Yiwei Li, Junhao Chen, Huawen Hu, Yihen Liu, Huaqin Zhao, Shaochen Xu, Haixing Dai, Lin Zhao, Ruidong Zhang, Wei Zhao, Zhenyuan Yang, Jingyuan Chen, Peilong Wang, Wei Ruan, Hui Wang, Huan Zhao, Jing Zhang, Yiming Ren, Shihuan Qin, Tong Chen, Jiaxi Li, Arif Hassan Zidan, Afrar Jahin, Minheng Chen, Sichen Xia, Jason Holmes, Yan Zhuang, Jiaqi Wang, Bochen Xu, Weiran Xia, Jichao Yu, Kaibo Tang, Yaxuan Yang, Bolun Sun, Tao Yang, Guoyu Lu, Xianqiao Wang, Lilong Chai, He Li, Jin Lu, Lichao Sun, Xin Zhang, Bao Ge, Xintao Hu, Lian Zhang, Hua Zhou, Lu Zhang, Shu Zhang, Ninghao Liu, Bei Jiang, Linglong Kong, Zhen Xiang, Yudan Ren, Jun Liu, Xi Jiang, Yu Bao, Wei Zhang, Xiang Li, Gang Li, Wei Liu, Dinggang Shen, Andrea Sikora, Xiaoming Zhai, Dajiang Zhu, and Tianming Liu. 2024. Evaluation of openai o1: Opportunities and challenges of AGI. CoRR, abs/2409.18486.

A Rationale and Details of the EM Algorithm

We provide the motivation and the details of applying the EM algorithm to optimize our framework.

A.1 Overview

We have two modules:

- **ReAligner**, parameterized by ψ , which generates Graph-aware Reasoning Chains z given (\mathcal{G},q) .
- Responser, parameterized by w, which predicts the final answer a given the question q and candidate Graph-aware Reasoning Chain z.

Given a training set of triples $\{(\mathcal{G}, q, a)\}$, our objective is to maximize:

$$\mathcal{O}(w, \psi) = \sum_{(\mathcal{G}, q, a)} \log \Big(\mathbb{E}_{z \sim p_{\psi}(z \mid \mathcal{G}, q)} \big[p_w(a \mid q, z) \big] \Big).$$

Because exact marginalization over z can be expensive, we employ an EM-style approach to iteratively refine both modules.

A.2 Rationale

The selection of the EM algorithm is fundamentally motivated by the central challenge and objective of the RAR: generating latent natural language (NL) reasoning steps for KGQA.

A.2.1 Core Challenge: Generating Latent NL Reasoning

- RAR aims to produce complex, human-like NL reasoning chains – the intermediate "thought process" connecting a question to its answer using a Knowledge Graph (KG).
- Crucially, these detailed reasoning steps are not available in standard KGQA training datasets. They constitute latent variables within our model.

A.2.2 Inadequacy of Direct Supervision Methods

• In standard SFT applied within RoG (Luo et al., 2024a), relation sequences absent from the original training data are often generated. To label these sequences for training, the process frequently relies on heuristics to create pseudo-gold labels, such as identifying the shortest KG path between query and answer entities, which can be potentially noisy.

 This heuristic-based supervision is insufficient for training models to generate the complex, multi-step, logically nuanced NL reasoning that RAR targets, especially when multiple constraints are involved.

A.2.3 EM as the Principled Approach for Latent Variables

- EM is the standard, principled approach for parameter estimation in models with latent variables.
- It provides a formal framework to optimize the Reasoner (generating the latent NL chain) and the Aligner (grounding the chain to the KG) without requiring explicit supervision for the intermediate NL reasoning steps.
- Optimization is guided indirectly by maximizing the likelihood of observing the correct final answer, conditioned on the feasibility of the generated reasoning chain being aligned to the KG.

A.2.4 Necessity of Iterative Refinement

- Generating coherent, long-form NL reasoning is challenging. Initial attempts, especially early in training, are likely to be imperfect or logically flawed.
- The iterative nature of the EM algorithm is well-suited for this progressive refinement: Estep identifies the most likely or "best" latent reasoning chains produced by the current model that successfully link the question to the correct answer via a feasible KG alignment. This step essentially evaluates the current reasoning quality based on outcomes; Mstep updates the parameters of the Reasoner and Aligner models by training them on these high-quality reasoning chains identified in the E-step. This step aims to make the models generate more chains similar to the successful ones.
- This iterative E-M loop allows the system to gradually improve the quality, logical coherence, and KG-alignability of the generated latent reasoning, as demonstrated qualitatively in Fig. 4.

A.2.5 Connections to Reinforcement Learning Connection to Implicit RL. The EM algorithm, as applied in RAR, can be viewed as a form of implicit Reinforcement Learning:

- The **E-step** acts like a selection or filtering mechanism based on the quality of the reasoning chain, implicitly assigning a high reward (*e.g*, 1) to successful chains (reaching the correct answer via KG alignment) and low reward (*e.g*, 0) to unsuccessful ones.
- The **M-step**, particularly the Reasoner update (maximizing $logp(z_hat_I|q)$ for selected high-quality chains z_hat_I), mathematically resembles a policy gradient update $(E[R*\nabla \log p(z|q)] \approx R*\nabla \log p(\hat{z}|q))$ where R is effectively this implicit binary reward.

Thus, EM reinforces the generation of "good" reasoning chains without the need for explicit reward engineering.

Why EM Was Preferred Over Explicit Reinforcement Learning. While the EM process here shares similarities with RL (see next point), we opted for EM over explicit RL formulations (like PPO) for several practical reasons:

- Reward Function Design: Crafting a good reward function ('R') that accurately captures the quality of multi-step NL reasoning is nontrivial.
- Training Complexity and Cost: Explicit RL methods often lead to higher computational costs and potentially unstable training.
- Efficiency and Simplicity: EM, derived naturally from the maximum likelihood objective for latent variable models, offers a more direct, mathematically grounded, and often simpler optimization pathway for our specific problem structure.

A.3 Details of EM Algorithm

A.3.1 Step 1: Updating Responser

1. Sample Candidate Graph-aware Reasoning Chains. For each training example (\mathcal{G}, q, a) , sample K Graph-aware Reasoning Chains:

$$z_k \sim p_\psi(z\mid \mathcal{G},q), \quad k=1,\ldots,K.$$
 Let $\hat{z}=\{\,z_1,z_2,\ldots,z_K\}.$

2. Approximate the Objective for w. The term

$$\log \mathbb{E}_{z \sim p_{\psi}(z \mid \mathcal{G}, q)} [p_w(a \mid q, z)]$$

is approximated by

$$\log \left(\frac{1}{K} \sum_{k=1}^{K} p_w(a \mid q, z_k)\right).$$

We then take gradients (w.r.t. w) and update w so that $p_w(a \mid q, z)$ is more likely to produce the correct a for the sampled Graph-aware Reasoning Chains.

3. **Result.** After updating w, Responser $p_w(a \mid q, z)$ is better aligned with whatever Graphaware Reasoning Chains p_{ψ} currently emits.

A.3.2 Step 2: EM-Style Update for ReAligner

After w is updated, we refine the ReAligner $p_{\psi}(z \mid \mathcal{G}, q)$. In EM terms, we view z as a latent variable: **E-Step (Posterior Inference)**

• Compute / Re-rank Graph-aware Reasoning Chains. Re-sample or re-rank the K Graph-aware Reasoning Chains using the updated p_w . We want Graph-aware Reasoning Chains that are "most aligned" with the correct answer a. Formally:

$$p_{w,\psi}(z \mid \mathcal{G}, q, a) \propto p_w(a \mid q, z) p_{\psi}(z \mid \mathcal{G}, q).$$

• **Scoring.** For a single Graph-aware Reasoning Chain *z*, define

$$S(z) = \log p_w(a \mid q, z) + \log p_{\psi}(z \mid \mathcal{G}, q).$$

• Selecting High-Quality Graph-aware Reasoning Chains. Rank (or sample) Graph-aware Reasoning Chains by S(z) and select the top set $z_I = \{\text{top-}K \text{ Graph-aware Reasoning Chains}\}.$

M-Step (Update ψ)

- Treat z_I (the selected high-quality Graphaware Reasoning Chains) as if they were observed.
- Update ψ by maximizing:

$$\log p_{\psi}(z_I \mid \mathcal{G}, q) = \sum_{z \in z_I} \log p_{\psi}(z \mid \mathcal{G}, q)$$

• In practice, this amounts to standard fine-tuning (e.g., instruction tuning or teacher forcing) of p_{ψ} on the newly identified high-quality Graph-aware Reasoning Chains.

A.3.3 Complete Iteration

A single iteration of our EM-style algorithm proceeds as follows:

1. (**Update** w): For each sample (\mathcal{G}, q, a) , draw K Graph-aware Reasoning Chains from p_{ψ} , then update w by maximizing

$$\log \left(\frac{1}{K} \sum_{k=1}^{K} p_w(a \mid q, z_k) \right).$$

2. (**E-Step**): Using the updated w, compute

$$p_{w,\psi}(z \mid \mathcal{G}, q, a) \propto p_w(a \mid q, z) p_{\psi}(z \mid \mathcal{G}, q).$$

Select high-quality Graph-aware Reasoning Chains z_I from the candidates.

3. (M-Step): Update ψ by maximizing $\log p_{\psi}(z_I \mid \mathcal{G}, q)$, i.e. fine-tune p_{ψ} so that it is more likely to emit z_I in the future.

This loop can be repeated until convergence or for a fixed number of epochs.

A.3.4 Practical Variations

- 1. **Top-**K **vs. Full Posterior.** Instead of summing/sampling over all subsets, it is simpler to pick the top-K Graph-aware Reasoning Chains by $S(\cdot)$.
- 2. Skipping Responser Optimization. To further improve efficiency, we can skip optimizing Responser. LLMs often possess strong zero-shot summarization or question-answering capabilities, which means they can produce high-quality answers from given Graph-aware Reasoning Chains without additional training. As a result, we can treat an LLM as a pre-optimized Responser and focus solely on updating ReAligner, thereby reducing overall computation.

B More Related Work

B.1 Comparison with Agent Exploration Methods

To situate RAR within the KGQA landscape, we first contrast it with representative agent exploration methods such as ToG (Sun et al., 2024). Although both paradigms comprise stages that can be informally mapped to *reasoning* and *grounding*, their internal principles diverge markedly.

Training methodology and optimization. RAR is trained with a mathematically grounded expectation—maximisation (EM) procedure that *explicitly and stably* refines two complementary capabilities: the NL *Reasoner* and the KG-aware *Aligner*. By contrast, many agent methods rely more heavily on prompt or workflow engineering (Cheng et al., 2024; Gu et al., 2024; Huang et al., 2024; Li et al., 2023; Nie et al., 2024); they seldom perform task-specific optimization that directly targets the core reasoning mechanism.

Accuracy, reliability, and complexity handling. The synergy between NL reasoning, KG-constrained alignment, and EM-guided supervised fine-tuning translates into markedly higher accuracy and robustness for RAR (see Tab. 1). Empirically, its explicit decomposition allows it to cope well with multi-hop and conjunctive constraints that are challenging for purely prompt-driven agents.

Resource consumption. Once training is finished, inference in RAR can be carried out by a collection of relatively small, specialized, fine-tuned models—one each for the Reasoner, Aligner, and Responser. This modularity yields the efficiency gains reported in Tab. ??. Agent systems, in contrast, often incur higher latency and cost because they perform several large-LLM calls while exploring the KG.

Taken together, these differences show that RAR is not a mere variant of the agent paradigm; its EM-centered optimization strategy and KG-constrained decoding constitute a distinct design that offers both practical efficiency and stronger empirical performance.

B.2 Comparison with Path Generation Methods

RAR also differs fundamentally from path generation methods such as RoG (Luo et al., 2024a) and GCR (Luo et al., 2024c). Where those systems directly predict a linear sequence of KG relations, ours maintains a higher-level, human-readable plan in natural language and lets the Aligner ground *each* step to KG triples.

Specifically, (i) the Reasoner produces a multi-step NL chain that expresses the conceptual logic; (ii) the Aligner incrementally matches every NL step to concrete triples through KG-constrained decoding; and (iii) the Responser integrates evidence from both the NL chain and the aligned KG

path to craft the final answer. Training again relies on EM—iteratively improving latent NL chains that can be aligned and that ultimately yield correct answers—whereas RoG and related work usually depend on direct SFT with shortest KG paths that may be noisy supervision signals.

Illustrative example. For the query "What did Dr Josef Mengele do?" the two paradigms unfold differently:

- *RoG*. An LLM planner outputs the relation people.person.profession; a symbolic retriever then follows this edge in the KG to obtain the triple (Josef Mengele, profession, Physician), and the answer *Physician* is returned.
- RAR. **Step 1** (A Reasoner step) proposes: "Identify the profession associated with *Josef Mengele*." The Aligner grounds this to the same triple as above. The Responser finally reports *Physician*, explicitly citing both the reasoning chain and the grounded path.

Handling complex constraints. Because RoG represents reasoning as a *single* linear relation sequence, it struggles with conjunctive queries such as "presidents who were also actors"; after following profession \rightarrow Actor, it cannot backtrack to verify profession \rightarrow President. RAR, in contrast, naturally decomposes the query into two successive NL steps ("find presidents" \rightarrow "filter those who are actors") and grounds each step separately, ensuring both constraints are satisfied.

Robustness to noisy supervision. Shortest-path supervision can be incorrect when more semantically plausible paths exist. By letting EM discover latent NL chains that are *both* alignable and answer-bearing, RAR avoids this brittleness and achieves the quality gains visualized in Fig. 4 of the submission.

In summary, the collaboration of an explicit NL Reasoner, a KG-constrained Aligner, and EM-based optimization endows RAR with a distinctive combination of interpretability, flexibility, and empirical strength that is not achieved by prior agent exploration or path generation methods.

C Experimental Setup

C.1 Fine-tuning Datasets and Knowledge Graph

For evaluation, we use two benchmark KGQA datasets: WebQuestionSP (WebQSP) (Yih et al., 2016) and Complex WebQuestions (CWQ) (Talmor and Berant, 2018). To ensure fair comparison, we adopt identical train and test splits as previous works (Jiang et al., 2022; Luo et al., 2024b). The detailed statistics of these datasets are presented in Tab. 6. Both WebQSP and CWQ are based on Freebase (Bollacker et al., 2008). To reduce computational overhead and memory consumption, we utilize the same subgraphs as previous work (Luo et al., 2024b). Additionally, we preprocess Freebase by converting CVT (Compound Value Type) nodes, which represent n-ary relationships, into binary relationships by concatenating edge labels with "-" as the delimiter, following Li et al. (2024).

C.2 Datasets for Cold Starting

To initialize the training of Reasoner and Aligner, we leverage high-quality Reasoning Chains and Knowledge Paths derived from SPARQL queries in the WebQSP and CWQ datasets. This approach prevents the models from generating malformed outputs during early training stages.

The SPARQL queries in these datasets represent the gold-standard Reasoning Chain that human experts would use to solve questions using the KG. We decompose these SPARQL queries according to a predefined grammar, breaking them into atomic chunks that each represent a single reasoning step. For each chunk, we query Freebase to obtain the corresponding triples, then use GPT-40 to generate natural language reasoning chains based on these retrieved triples. Through this process, we generate a dataset of 2,000 high-quality Reasoning Chains with their corresponding Knowledge Paths for each question. This dataset enables us to perform coldstart pretraining of Reasoner and Aligner, teaching them to generate well-structured, step-by-step Reasoning Chains and Knowledge Paths with appropriate special tokens, a crucial foundation for subsequent optimization using the EM algorithm.

C.3 Hyperparameters

For Responser, we adopt Llama-3.1-8B (Meta, 2024) without fine-tuning based on our preliminary experiments (detailed analysis in App. A.3.4). For both Reasoner and Aligner, we conduct exten-

sive experiments with various lightweight LLMs ranging from 0.5B to 8B parameters (Yang et al., 2024; Touvron et al., 2023; Meta, 2024). All models share the same hyperparameter configuration: training for 3 epochs with a batch size of 4 and a learning rate of 2e-5. We employ a cosine learning rate scheduler with a warmup ratio of 0.03. The training is performed on 4 A6000 GPUs for each model variant.

D Details of KG-constrained Decoding

For efficient knowledge graph operations, we implemented a Virtuoso-based Freebase instance with distributed high-speed SPARQL querying capabilities. Our system achieves a throughput of 2,000 requests per second, enabling rapid graph traversal and neighborhood node retrieval during the constrained decoding process. This high-performance infrastructure allows us to efficiently retrieve nexthop candidates for token constraints directly from the knowledge graph.

E Case Study of Different Iteration Steps

In Fig. 7, we provide two examples to investigate the effect of iteration steps of the EM algorithm.

F Templates and Prompts

In this section, we illustrate all the prompts used in the experiments.

Reasoning Chain Template. The template of Reasoning Chains generated by Reasoner is shown in Fig. 8, where each s_i is a discrete reasoning step in natural language.

Knowledge Path Template. The template of Knowledge Path generated by Aligner is shown in Fig. 9, where e abd r denotes the entities and relations from the KG.

ReAligner Prompt. The prompt for instructing ReAligner is shown in Fig. 10, where the task is to generate the Reasoning Chain and Knowledge Path given the question and the KG.

Responser Prompt. The prompt for instructing Responser is shown in Fig. 11, where the task is to generate the final answer based on the given the question and the generated Reasoning Chain and Knowledge Path.

LLM-driven Consolidation Prompt. The prompt for LLM-driven Consolidation is shown in Fig. 12. We use RAR to generate K Knowledge Paths and hypothesis answers for each question. The Knowledge Paths and hypothesis answers are provided

Dataset	Dataset Statistics		Statistics of Answer Numbers			
	#Train	#Test	#Ans = 1	$2 \ge \#Ans \le 4$	$5 \ge \#Ans \le 9$	#Ans ≥ 10
WebQSP CWQ	2,826	1,628	51.2%	27.4%	8.3%	12.1%
CWQ	27,639	3,531	70.6%	19.4%	6%	4%

Table 6: Statistics of datasets.

Case 1	Question: What high school did the artist who recorded "Girl Tonight" attend? Answer: Petersburg.			
Steps	20	100		
Reasoning Chain & Knowledge Path	First, identify the artist "c" associated with the recording "Girl Tonight". Sohn H. Guyer, school type, High school Next, determine the high school attended by the artist "c". Sohn H. Guyer, notable_object, High school	First, identify the artist "c" associated with the recording"Girl Tonight". Arey Songz, recordings, Girl Tonight > Next, determine the educational institutions "x" attended by artist "c". Arey Songz, education_institution, Petersburg> Finally, confirm which of educational institutions "x" is a high school. Petersburg, school_type, High school>		
Response	John H. Guyer X	Petersburg ✓		
Case 2	Question: What type of Claude Debussy music appears in the film Black Tights? Answer: Ballet.			
Steps	20	100		
Reasoning Chain & Knowledge Path	First, identify the type of music "x" associated with Claude Debussy. Claude Debussy, compositions, En blanc et noir> Next, confirm that the type of music "x" is featured in the film Black Tights. 	First, identify the type of music "x" associated with Claude Debussy. Claude Debussy, genre, Ballet> Next, determine whether the type of music "x" appears in the film Black Tights. Sallet, films_in_this_genre, Black Tights>		
Response	En blanc et noir X	Ballet ✓		

Figure 7: Examples of Reasoning Chains and Knowledge Paths generated by RAR under different iteration steps.



Figure 8: The template for the Reasoning Chain generated by Reasoner.

to general LLMs to answer the questions without fine-tuning.

Knowledge Path Template

 $<\!\!\text{ALIGN}><\!\!\text{TRIPLE}><\!\!\text{I}>\!\!e_1<\!\!\text{I}>\!\!r_1<\!\!\text{I}>\!\!e_1'<\!\!\text{TRIPLE}>\ldots<\!\!\text{TRIPLE}><\!\!\text{I}>\!\!e_n<\!\!\text{I}>\!\!r_n<\!\!\text{I}>\!\!e_n'<\!\!\text{TRIPLE}><\!\!\text{ALIGN}>$

Figure 9: The template of Knowledge Paths generated by Aligner.

ReAligner Prompt
Generate a step-by-step thinking process for the given question. Ensure the thinking process is aligned with triples in the knowledge base.
Question: <question></question>
Query entities: <question entities=""></question>
======================================
Align Process: <knowledge path=""></knowledge>

Figure 10: The prompt template for ReAligner.

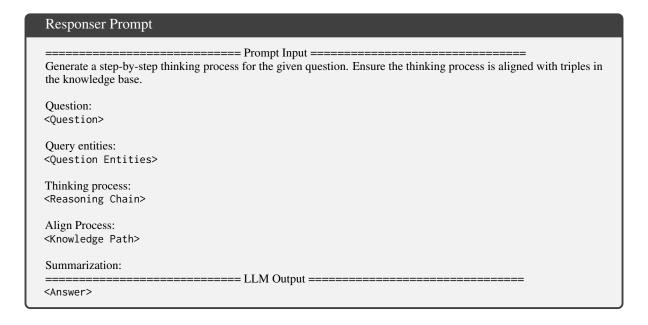


Figure 11: The prompt template for Responser.

Figure 12: The prompt template for LLM-driven Consolidation.