

Diagnosing and Remediating Representation Deficiencies for Deterministic Reasoning in KGQA

Anonymous ACL submission

Abstract

Large language models (LLMs) have demonstrated increasingly strong reasoning capabilities, achieving remarkable progress in knowledge graph question answering (KGQA). However, a key challenge in such systems is non-deterministic reasoning, where the model indecisively activates multiple semantically related knowledge graph edges for a given query, frequently leading to incorrect answers. To address this issue, we propose **Diagnosing and Remediating Representation Deficiencies for Deterministic Reasoning in KGQA (DR²)**. DR² identifies and localizes non-deterministic reasoning behaviors, uncovering the underlying semantic representation deficiencies in LLMs. Building on this diagnosis, we design abductive reasoning-based preference learning, which promotes fine-grained semantic discrimination and mitigates non-deterministic reasoning errors. Experimental results demonstrate that the proposed DR² significantly outperforms several strong baselines, achieving state-of-the-art performance on the widely used WebQSP and CWQ benchmarks.¹

1 Introduction

Large language models have demonstrated increasingly strong reasoning capabilities, leading to remarkable improvements in Knowledge Graph Question Answering (KGQA) tasks. Current approaches primarily leverage LLMs’ powerful capabilities through two paradigms: in-context learning augmented with retrieved evidence to directly generate answers or reasoning chains (Li et al., 2023; Nie et al., 2024), and interactive, step-by-step frameworks where the model decomposes questions or traverses the knowledge graph through tool calls (Gu et al., 2023; Huang et al., 2024; Sun et al., 2024). To further enhance precision, supervised fine-tuning on questions and

¹Our codes and data will be released upon acceptance.

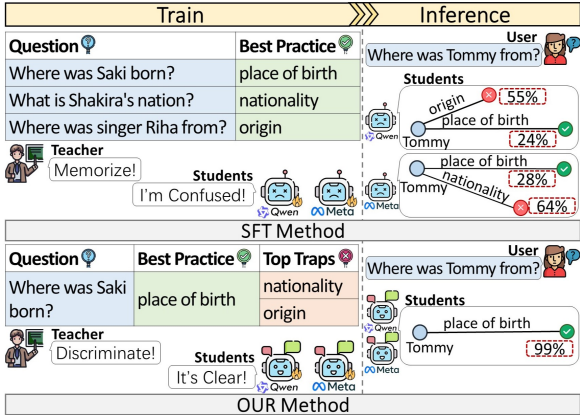


Figure 1: Comparing SFT’s memorization-induced indecisiveness with our method’s deterministic discrimination on semantically related edges.

synthesized reasoning paths is widely used to enable more accurate and direct generation of search queries (LUO et al., 2024; Luo et al., 2024; Bu et al., 2025; Xu et al., 2025b). These methods advance KGQA by more effectively combining the knowledge of LLMs with the knowledge graphs, enhancing the reliability of reasoning process.

However, a key challenge in such systems is non-deterministic reasoning, where the model indecisively activates multiple semantically related knowledge graph edges for a given query, frequently leading to incorrect answers. Existing Supervised Fine-Tuning (SFT) approaches in KGQA fundamentally fall short of addressing this semantic ambiguity. As Figure 1 illustrates, the indecisiveness of SOTA models on semantically similar knowledge graph edges—such as confusing “place of birth” with “nationality”—accounts for up to 53% of overall errors (see Figure 7). This critical failure stems from a key limitation of SFT: while it optimizes for matching the correct tool call during training, it does not equip the model with the ability to internally discriminate between closely related relations. We empirically

065 verify this in a controlled setting with the easily
066 confusable tool calls “language spoken” and
067 “official language”. Although SFT method achieves
068 high accuracy on the training set, the performance
069 on unseen test data drops to 75%, as shown in
070 Figure 8. This 25% accuracy gap reveals the
071 model’s reliance on memorization rather than
072 true semantic discrimination, and its persistent
073 vulnerability to non-deterministic reasoning when
074 encountering unfamiliar examples.

075 To mitigate this issue, we propose Diagnosing
076 and Remediating Representation Deficiencies for
077 Deterministic Reasoning in KGQA (**DR**²), which
078 enhances the model’s ability to discriminate
079 between semantically related knowledge graph
080 edges. **DR**² identifies and localizes non-
081 deterministic reasoning behaviors, uncovering the
082 underlying semantic representation deficiencies
083 in LLMs. Building on this diagnosis, we
084 design abductive reasoning-based preference
085 learning, which promotes fine-grained semantic
086 discrimination and mitigates non-deterministic
087 reasoning errors. **DR**² enhances the model’s
088 capability to distinguish fine-grained semantic
089 differences, our main contributions are:

- 090 • We propose a novel method **DR**² to detect non-
091 deterministic reasoning and locate semantic
092 representation deficiencies in LLMs.
- 093 • We propose using abductive reasoning method
094 to acquire fine-grained semantic preference
095 data to mitigate the deficiencies.
- 096 • The proposed **DR**² method achieved state-
097 of-the-art performance on two widely used
098 benchmarks WebQSP and CWQ.

099 2 Related Works

100 **Preferences Alignment.** Recent advances in
101 LLM alignment that leverage human preferences
102 offer a promising direction for enhancing the
103 discriminative capability of KBQA systems.
104 The established method Reinforcement Learning
105 from Human Feedback (RLHF) trains a reward
106 model to optimize the policy model (Ouyang
107 et al., 2022; Bai et al., 2022). Direct
108 Preference Optimization (DPO) eliminates the
109 separate reward model by training directly on
110 output preference pairs (Rafailov et al., 2023).
111 Subsequent work simplified DPO by removing
112 the reference model (Meng et al., 2024; Hong
113 et al., 2024). Most recently, IOPO improves

114 DPO by constructing bidirectional preference pairs,
115 improving the model’s fine-grained discriminative
116 capability (Zhang et al., 2025).

117 **Knowledge Graph Question Answering.** Early
118 research can be categorized into rule-based
119 methods, that leverage the symbolic structure of
120 Knowledge Bases(KB) for retrieval or parsing (Sun
121 et al., 2018; Zhang et al., 2022; Ye et al.,
122 2022) and embedding-based methods which utilize
123 neural networks to learn latent representations
124 of entities and relations for reasoning (Miller
125 et al., 2016; Yasunaga et al., 2021; Jiang et al.,
126 2022). With the advent of large language
127 models (LLMs), KBQA methods have evolved
128 significantly. Current approaches predominantly
129 leverage LLMs’ powerful reasoning capabilities,
130 primarily through in-context learning (ICL)
131 augmented with retrieved evidence to generate
132 answers or reasoning chains directly (Li et al.,
133 2023; Nie et al., 2024), as well as through step-
134 by-step frameworks where the model interactively
135 decomposes questions or traverses the knowledge
136 graph (Gu et al., 2023; Huang et al., 2024; Sun
137 et al., 2024). To further enhance performance and
138 alignment with specific KB schemas, supervised
139 fine-tuning (SFT) on datasets of question-logical
140 form pairs has been effectively adopted to enable
141 more accurate and direct generation of structured
142 queries (Luo et al., 2024; LUO et al., 2024; Jiang
143 et al., 2025; Xu et al., 2025b; Bu et al., 2025).

144 While fine-tuning enhances the model’s ability
145 to understand questions and produce formally
146 correct queries, the fine-tuned models may still lack
147 robust discriminative capabilities, often leading
148 to non-deterministic reasoning when faced with
149 semantically similar edges. To mitigate this
150 issue, we propose the Diagnosing and Remediating
151 Representation Deficiencies for Deterministic
152 Reasoning in KGQA to explicitly enhance the
153 model’s ability to discriminate between edges with
154 close semantic meanings, thereby improving the
155 precision and reliability in KBQA tasks.

156 3 Preliminary

157 In this section, we introduce several basic concepts
158 and definitions used in our work.

159 **Knowledge Graph(KG).** A Knowledge Graph
160 is a set of triples, denoted as $KG =$
161 $\{(e_0, r, e_1) | e_0, e_1 \in \mathcal{E}, r \in \mathcal{R}\}$, where \mathcal{E} is the
162 set of entities and \mathcal{R} is the set of relations.

163 **Reasoning Path.** Given a query, a reasoning path

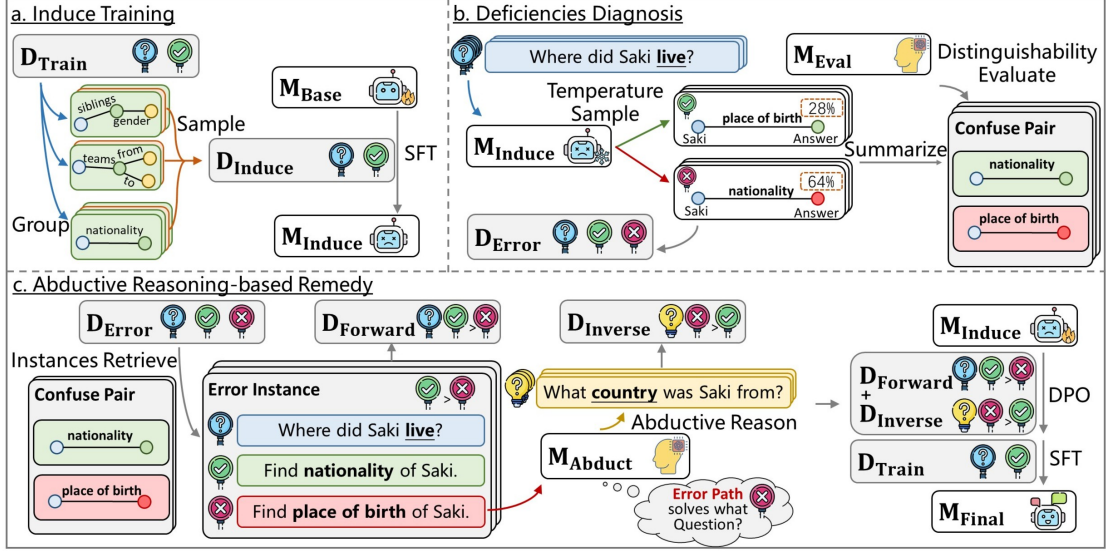


Figure 2: The overall framework of our DR². The diagnosing process trains an induced model with limited data (a) and diagnoses its deficiencies by analyzing the reasoning behaviors under temperature sampling (b). The remedying process leverages abductive reasoning to construct preference data for training the final model (c).

P is a sequence from the topic entity to the answer entity. It is formally represented as $P = e_{topic} \xrightarrow{r_1} e_1 \xrightarrow{r_2} \dots \xrightarrow{r_n} e_{ans}$, where e_{topic} is the topic entity and e_{ans} is the answer entity. In the training set D_{Train} , the golden reasoning path P_{gold} is directly derived from its annotated SPARQL query.

Reasoning Structure. The reasoning structure is an abstraction of a reasoning path, obtained by removing all entity nodes and retaining only the sequence of relations. For a path $P = e_{topic} \xrightarrow{r_1} e_1 \xrightarrow{r_2} \dots \xrightarrow{r_l} e_{ans}$, its reasoning structure is defined as $P^S = [r_1, r_2, \dots, r_l]$.

4 Approach

In this section, we provide a detailed description of our proposed DR² method. As illustrated in Figure 2, our approach to mitigating semantic representation deficiencies is decomposed into three tasks: Inducing Training, Deficiencies Diagnosis and Abductive Reasoning-based Remedy.

4.1 Inducing Training

Inducing Training aims to train an "induced model" to systematically expose the model's inherent semantic representation deficiencies. As shown in Figure 8, the model exhibits non-deterministic reasoning during early SFT, where the probabilities for the two similar paths rise synchronously. To induce this behavior, we design a training strategy based on reasoning structure classification.

Reasoning Structure Classification. For each

data (Q, P_{Gold}) in the training set D_{Train} , we compute its reasoning structure P_{Gold}^S . We then partition D_{Train} into K mutually exclusive subsets $\{G_1, G_2, \dots, G_K\}$, where all data within the same subset share the same reasoning structure.

Non-determinism Inducing Training. To ensure the model is exposed to all reasoning structures, we randomly sample a small proportion of data from each subset G_i according to a pre-defined ratio to form the inducing training set D_{Induce} . We then fine-tune the base model M_{Base} using D_{Induce} to obtain the induced model M_{Induce} .

4.2 Deficiencies Diagnosis

After obtaining the induced model M_{Induce} , we are able to diagnose the model's semantic representation deficiencies by analyzing its non-deterministic reasoning. We first collect confusion errors from M_{Induce} , and summarize them into confusing structure pairs. A Semantic Distinguishability Evaluation is then proposed to pinpoint the pairs stemming from inherent semantic representation deficiencies, providing precise targets for remediation.

Confusion Error Collection. For each data (Q, P_{Gold}) in the training set D_{Train} , we feed the question Q into the induced model M_{Induce} and obtain m reasoning paths with a temperature sampling ($Sample_{Temp}$) process:

$$P_{Gen} = Sample_{Temp}(M_{Induce}, Q). \quad (1)$$

If the structure of the generated path differs from

that of the golden path, it forms a confusion error instance (Q, P_{Gold}, P_{Gen}) . All such instances constitute the confusion error instance set D_{Error} .

To focus on systematic errors caused by semantic representation deficiencies, we summarize errors at the reasoning structure level. For each error instance (Q, P_{Gold}, P_{Gen}) in D_{Error} , we extract the reasoning structure of the golden path $S_1 = P_{Gold}^S$ and the generated path $S_2 = P_{Gen}^S$, then (S_1, S_2) forms an unordered confusion pair.

Semantic Distinguishability Evaluation. To pinpoint the confusion that arises from intrinsic deficiencies of the model, we evaluate the semantic distinguishability of each summarized confusion pair (S_1, S_2) . We first construct two question groups G_1 and G_2 which contains all questions from D_{Train} whose golden reasoning structure is included in S_1 and S_2 , respectively:

$$G_i = \{Q | (Q, P) \in D_{Train} \wedge P^S = S_i\}. \quad (2)$$

We randomly sample a question Q from $G_1 \cup G_2$, and then classify it into G_1 or G_2 using a general-purpose LLM M_{Eval} . The average classification accuracy (Acc) over multiple sampled questions yields the distinguishability score $DistSc$ using:

$$DistSc = Acc(M_{Eval}(Q) = Group(Q)). \quad (3)$$

A high $DistSc$ score indicates that the questions are semantically distinct, and the model’s confusion stems from its representation deficiency. Pairs with scores above a predefined threshold form the target deficiencies set T_{Defect} , which is used for further remediation.

4.3 Abductive Reasoning-based Remedy

For each target deficiency in T_{Defect} , we extract the original error instances from D_{Error} as the forward preference data. We then leverage the abductive reasoning capability of a general-purpose LLM to construct inverse preference data. Finally, joint training on both forward and inverse preferences guides the model to deeply understand and discriminate between the confusion pairs, thereby remedying its representation deficiencies.

Error Instances Retrieval. For each target confusion pair in T_{Defect} , we filter all matching error instances (Q, P_{Gold}, P_{Gen}) from the error set D_{Error} . For each instance, a forward preference tuple $(Q, P_{Gold} \succ P_{Gen})$ is constructed, which is added to the forward preference set $D_{Forward}$.

Inverse Question Generation. For each matched error instance (Q, P_{Gold}, P_{Gen}) , we employ a general-purpose LLM M_{Abduct} to perform abductive reasoning, inferring a new natural language question Q_{Inv} for which P_{Gen} becomes the uniquely correct reasoning path:

$$Q_{Inv} = M_{Abduct}(P_{Gen}). \quad (4)$$

This process yields a inverse preference tuple $(Q_{Inv}, P_{Gen} \succ P_{Gold})$, which is added to the inverse preference set $D_{Inverse}$.

Preference Learning. The induced model M_{Induce} is optimized using Direct Preference Optimization on the union of the forward preference set $D_{Forward}$ and the inverse preference set $D_{Inverse}$. This directly targets the identified semantic representation deficiencies. The model is then Supervised Fine-Tuned on the original training set D_{Train} to produce the final model M_{Final} , which enhances its overall reasoning ability.

5 Experiment

In this section, we present a comprehensive experimental evaluation of the proposed DR² method. We first introduce the experimental setup, including the benchmarks, evaluation metrics, and baseline methods. Following this, we report the main results and a series of analytical experiments to examine the characteristics of DR² against baseline methods from multiple perspectives.

5.1 Experimental Settings

Benchmarks. To evaluate the performance of our proposed method on KGQA tasks, we employ two widely used benchmarks: WebQSP (Yih et al., 2016) and CWQ (Talmor and Berant, 2018).

Metrics. We adopt standard metrics for KGQA evaluation: Hits, F1 and Hits@1 following previous works. In our work, Hits measures if any predicted answer is correct, while Hits@1 specifically checks if the first predicted answer is correct. See detailed definitions in Appendix C.

Baselines. We compare DR² against previous state-of-the-art and representative KGQA methods, including strong baselines such as MemQ (Xu et al., 2025b) and KaeDe (Bu et al., 2025), as well as other notable methods for broader context.

Implementation Details. To ensure a fair comparison, we use Llama2-7B-chat (Touvron et al., 2023) as our base language model, consistent with the setup in recent methods like MemQ (Xu

Method	Society		Entertainment		Culture		Average	
	Hits@1	F1	Hits@1	F1	Hits@1	F1	Hits@1	F1
KaeDe	0.620	0.632	0.519	0.634	0.696	0.663	0.589	0.640
MemQ	0.741	0.773	0.561	0.648	0.696	0.790	0.648	0.718
DR ²	0.759	0.808	0.603	0.677	0.732	0.798	0.681	0.745

Table 1: Performance comparison of DR² and baselines on a confusing subset of WebQSP and CWQ.

et al., 2025b) and KaeDe (Bu et al., 2025). All our experiments are conducted on a machine with 4 NVIDIA RTX 4090 GPUs. During inference, we use a beam size of 4 for both datasets.

5.2 Main Result

Method	WebQSP		CWQ	
	Hits	F1	Hits	F1
KV-Mem (Miller et al., 2016)	0.467	0.345	0.184	0.157
GraftNet (Sun et al., 2018)	0.664	0.604	0.368	0.327
QGG (Lan and Jiang, 2020)	0.730	0.738	0.369	0.374
DECAF (Yu et al., 2022)	0.821	0.788	-	-
UniKGQA (Jiang et al., 2022)	0.751	0.702	0.507	0.480
ToG (Sun et al., 2024)	0.826	-	0.676	-
KG-Agent (Jiang et al., 2025)	0.833	0.810	0.722	0.692
FiDeLiS (Sui et al., 2025)	0.844	0.783	0.715	0.643
AMAR (Xu et al., 2025a)	0.843	0.812	0.831	0.785
FM-KBQA (Gao et al., 2025)	0.873	0.842	0.795	0.687
RoG (LUO et al., 2024)	0.857	0.708	0.626	0.562
ChatKBQA (Luo et al., 2024)	0.864	0.835	0.860	0.813
GGI-MAB (Tang et al., 2025)	0.866	0.756	0.794	0.720
EPERM (Long et al., 2025)	0.888	0.724	0.662	0.589
KaeDe (Bu et al., 2025)*	0.908	0.819	0.880	0.801
MemQ (Xu et al., 2025b)*	0.890	0.860	0.878	0.836
DR ²	0.918	0.889	0.898	0.860

Table 2: The results of our method DR² compared with previous approaches on WebQSP and CWQ. The asterisk * denotes the results we reproduced.

The main results of our proposed DR² on the WebQSP and CWQ datasets benchmarks are presented in Table 2 and Table 3. Our method achieves state-of-the-art performance across all key metrics on both benchmarks. This consistent and superior performance provides empirical evidence that our approach effectively mitigates the semantic representation deficiency in LLMs, confirming the effectiveness of diagnosing and remedying representation deficiencies in enhancing deterministic reasoning for KGQA. These results not only demonstrate the practical advantages of our method, but also highlight the fundamental importance of semantic representation robustness in complex reasoning tasks with LLMs.

Method	WebQSP	CWQ
RoG (LUO et al., 2024)	0.795	0.567
KaeDe (Bu et al., 2025)	0.830	0.798
MemQ (Xu et al., 2025b)	0.847	0.805
DR ²	0.882	0.835

Table 3: The results of our method DR² compared with previous approaches on WebQSP and CWQ using the strict Hits@1 metric (see Appendix C) for details.

5.3 Discrimination Capability Evaluation

To further evaluate the capability of DR² in discriminating between semantically related knowledge graph edges, we constructed a challenging subset derived from the WebQSP and CWQ test sets. This subset comprises 671 questions whose gold reasoning paths contain knowledge graph edges that are hard to distinguish. The questions are categorized into three domains (Society, Entertainment and Culture), based on their topic entities.

As shown in Table 1, DR² achieves consistent improvements over all baseline methods across all evaluation metrics on the challenging subset. These results confirm that our method effectively mitigates the semantic representation deficiency for semantically related edges.

5.4 Mitigation of Representation Deficiencies

To provide direct evidence of how DR² mitigates representation deficiencies, we compare its internal representations with those of MemQ on the identified defect set T_{Defect} . For each diagnosed defect pair (S_1, S_2) , we input the corresponding questions and golden paths into the model to extract the hidden state of the last token, which forms representation groups G_1 and G_2 . We then project the representations onto the first two principal components using Principal Component Analysis and visualize the distributions with Kernel Density Estimation. In addition, we calculate the Separability Score ($SepSc$) to quantify this difference, defined as the ratio of average inter-

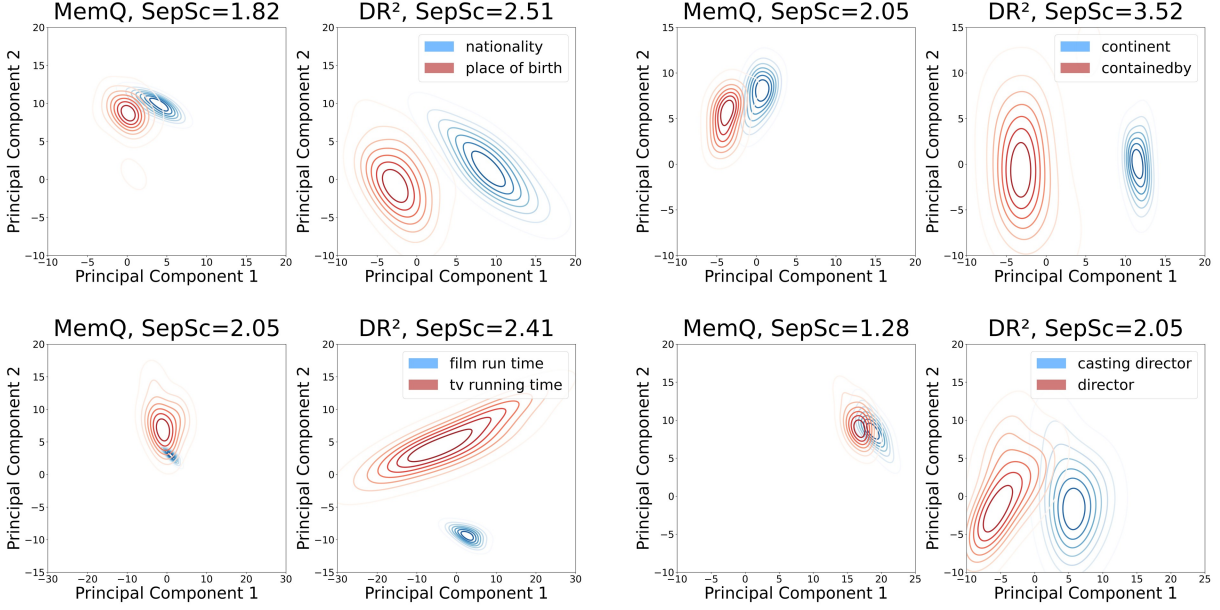


Figure 3: Comparing the separability of representations between DR² and MemQ.

group Euclidean distance (G) to average intra-group Euclidean distance (D):

$$SepSc(G_1, G_2) = \frac{GD(G_1, G_2)}{\frac{1}{2} \times (D(G_1) + D(G_2))}. \quad (5)$$

As shown in Figure 3, compared to MemQ, DR² produces more separated representation clusters for the structures within each defect pair compared to MemQ, with a clearer boundary between them. Quantitatively, DR² achieves a consistently higher $SepSc$ for all these pairs. This demonstrates that our method effectively enhances the model’s ability to discriminate between the confusable reasoning structures. See further analyses in Appendix E.

5.5 Enhancement of Deterministic Reasoning

To quantitatively assess the improvement in reasoning determinism, we construct a challenge test set comprising questions that are prone to cause confusion. We analyze the model behaviors on this set by performing temperature sampling for both MemQ and DR². A model is considered to perform deterministic reasoning if it generates a path with a probability greater than 90%.

The histogram in Figure 4 presents the distribution of probabilities that the model generates the correct reasoning path. The result shows that DR² increases the probability of deterministically inferring the correct path by 65% over MemQ. This confirms that our method effectively enables the model to perform deterministic reasoning towards the correct path.

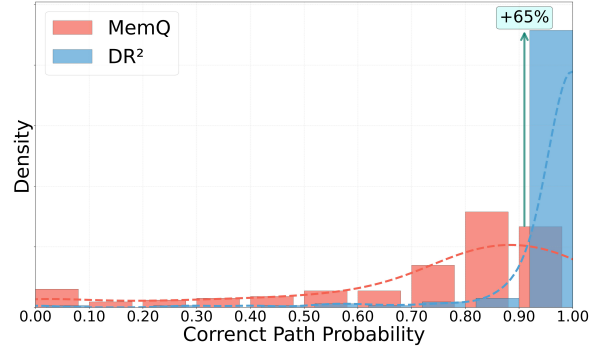


Figure 4: Comparing reasoning determinism of MemQ and DR² via temperature sampling on a confusion set.

5.6 Cross-Model Effectiveness

To assess the generalizability of DR² beyond a single model, we evaluate its performance when using different base models. We conduct experiments on four widely-used LLMs, the results in Table 4 demonstrate that all evaluated models achieve strong performance. This indicates that our approach to diagnosing and remedying semantic representation deficiencies is not restricted to a single model architecture, but instead exhibits effective transferability across a variety of backbone models. These findings further establish DR² as a robust and versatile framework for enhancing deterministic reasoning in KGQA, regardless of the underlying language model.

5.7 Ablation Study

To validate the contribution of each component in our proposed DR² method, we conduct a

Base	Method	WebQSP			CWQ		
		Hits@1	F1	Hits	Hits@1	F1	Hits
Llama2	MemQ	0.847	0.860	0.890	0.805	0.836	0.878
	DR ²	0.882 (+0.035)	0.889 (+0.029)	0.918 (+0.028)	0.835 (+0.030)	0.860 (+0.024)	0.898 (+0.020)
Llama3	MemQ	0.868	0.874	0.895	0.814	0.839	0.879
	DR ²	0.882 (+0.014)	0.888 (+0.014)	0.919 (+0.024)	0.840 (+0.026)	0.868 (+0.029)	0.911 (+0.032)
Qwen	MemQ	0.830	0.845	0.877	0.796	0.822	0.868
	DR ²	0.840 (+0.010)	0.857 (+0.012)	0.899 (+0.022)	0.808 (+0.012)	0.837 (+0.015)	0.889 (+0.021)
Gemma	MemQ	0.854	0.856	0.879	0.799	0.826	0.867
	DR ²	0.870 (+0.016)	0.876 (+0.020)	0.905 (+0.026)	0.809 (+0.010)	0.839 (+0.013)	0.883 (+0.016)

Table 4: Evaluation of the generalization of DR² across multiple widely-used LLMs.

comprehensive ablation study with the following four experimental settings: **1) Without Inverse Question Generation (w/o IQG)**. During the preference construction stage, only the forward preference data $D_{forward}$ are used for DPO training, thereby isolating the contribution of the synthesized counterfactual data. **2) Without Semantic Distinguishability Evaluation (w/o SDE)**. It eliminates the semantic distinguishability evaluation in the diagnosis stage. All summarized structural confusion pairs are directly used as the target defect set T_{Defect} for subsequent correction. **3) Without Confusion Error Summarization (w/o CES)**. It uses all sampled confusion error instances from D_{Error} are directly used to construct forward preferences for DPO training. **4) Without Direct Preference Optimization (w/o DPO)**. The model is directly fine-tuned on the original training set D_{Train} using a SFT objective.

Strategy	WebQSP			CWQ		
	Hits@1	F1	Hits	Hits@1	F1	Hits
DR ²	0.882	0.889	0.918	0.835	0.860	0.898
w/o IQG	0.875	0.882	0.913	0.820	0.850	0.892
w/o SDE	0.864	0.873	0.905	0.819	0.850	0.889
w/o CES	0.844	0.861	0.904	0.812	0.845	0.887
w/o DPO	0.848	0.861	0.891	0.801	0.833	0.876

Table 5: Ablation study on the components of DR².

Based on the results in Table 5, we observe the following: 1) The performance gap between DR² and "w/o IQG" demonstrates that the generated inverse preference data are essential for capturing fine-grained semantic distinctions. 2) The further decline in performance from "w/o IQG" to "w/o SDE" highlights the necessity of semantic distinguishability evaluation; without this step, the refinement process is adversely affected by confusion pairs unrelated to inherent model deficiencies. 3) The comparison between

"w/o Semantic Evaluation" and "w/o Confusion Summarization" reveals that using all error instances without abstraction impedes the model's ability to focus on the most frequent and systematic confusion pairs, resulting in diminished performance. 4) Most notably, the largest performance drop occurs in the "w/o DPO" setting, confirming that standard SFT alone is inadequate for addressing the underlying semantic representation deficiencies in LLMs.

5.8 Reasoning Precision Improvement

Total Hops	1	2	3,4	>=5	avg
<i>EHR</i>					
KaeDe	0.716	0.770	0.759	0.679	0.741
MemQ	0.800	0.817	0.828	0.839	0.821
DR ²	0.875	0.867	0.866	0.869	0.868
<i>GoldGED</i>					
KaeDe	0.157	0.420	0.744	2.078	0.654
MemQ	0.170	0.358	0.659	1.118	0.523
DR ²	0.155	0.315	0.619	1.075	0.490

Table 6: Evaluation of edge accuracy (*EHR*) and search structure (*GoldGED*) for DR² and baselines.

To further investigate the improvement in the model's reasoning capability, we evaluate the reasoning paths generated by the model against the golden reasoning paths. The evaluation focuses on two key aspects: 1) the accuracy of the edges and 2) the accuracy of the search structure.

We use the Edge Hit Rate (*EHR*) to evaluate the accuracy of edges, which is defined as the proportion of edges in the golden path that are correctly predicted in the predicted path.

$$EHR = \frac{\text{num}(\{r|r \in P_{pred} \wedge r \in P_{gold}\})}{\text{num}(\{r|r \in P_{gold}\})}. \quad (6)$$

We use the Graph Edit Distance (*GoldGED*) to evaluate the accuracy of the search structure, which is defined as the edit distance between the predicted

path and the golden path. A lower *GoldGED* indicates a higher structural accuracy.

$$GoldGED = \min_{\pi \in \Pi(P_{pred}, P_{gold})} num(\pi). \quad (7)$$

The evaluation results are presented in Table 6. In the comparison of edge accuracy, DR² method achieves a significantly higher *EHR* than all baseline methods. Moreover, as the complexity of the questions increases, the *EHR* of our method remains at a comparable level without a noticeable decline, demonstrating its robustness. In the comparison of graph structure, our method attains a lower *GoldGED* than baselines, confirming that the search graphs generated by our approach predicts reasoning path structures more accurately.

5.9 Error Analysis

To gain deeper insight into the specific improvements introduced by our method, we conduct a detailed error analysis on both datasets. We categorize errors according to the source of reasoning confusion: **(1) Main Path Confusion**, where errors occur along the primary reasoning path from the starting entity to the answer entity; and **(2) Filter Path Confusion**, which arises from the misapplication of filtering constraints.

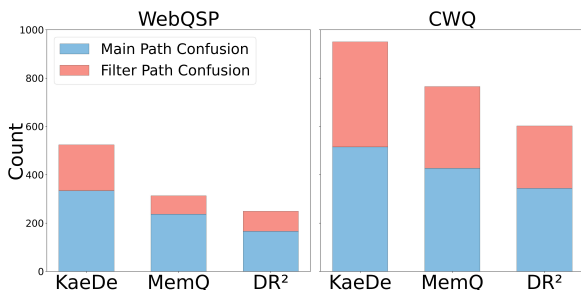


Figure 5: Error analysis of main path and filter path confusion for DR² and baselines.

As shown in Figure 5, DR² yields fewer errors in both Main Path and Filter Path confusion compared to the strong baselines MemQ and KaeDe. This demonstrates that the enhanced fine-grained semantic discrimination enabled by our method effectively reduces core relation misidentification and improves the precision of constraint application, thereby leading to more robust and reliable reasoning.

5.10 Case study

We present a case to intuitively demonstrate how DR² enhances deterministic reasoning. As shown

WebQTest-1370	
Question:	where did thomas hobbes live?
Golden Path:	Find where *Thomas Hobbes* has lived.
DR ²	
Generated Path: (<i>P</i> = 97.56%)	Find where *Thomas Hobbes* has lived.
MemQ	
Generated Path 1: (<i>P</i> = 51.92%)	Find place of birth of *Thomas Hobbes*.
Generated Path 2: (<i>P</i> = 44.53%)	Find where *Thomas Hobbes* has lived.

Figure 6: Case study comparing the deterministic reasoning of DR² and MemQ on a representative query.

in Figure 6, for the query “Where did Thomas Hobbes live?”, the golden reasoning path is “Find the location where Thomas Hobbes has lived”. The structural pair “the location where someone has lived” and “the place of birth of someone” was identified as a semantic representation deficiency during diagnosis. In this case, DR² demonstrates clear deterministic reasoning by assigning a high probability of 97.56% to the correct path, while MemQ exhibits pronounced indecision, assigning 51.92% to the incorrect path “Find the place of birth of the person Thomas Hobbes” and 44.53% to the correct one. See more cases in Appendix F.

6 Conclusion

In this paper, we propose DR², a novel method that mitigates the key challenge of non-deterministic reasoning in KGQA by diagnosing and remedying underlying representation deficiencies. Our approach diagnoses such deficiencies by analyzing the model’s non-deterministic behavior under temperature sampling. These identified deficiencies are remedied by constructing targeted preference data through abductive reasoning-based Preference Learning. Experimental results demonstrate that DR² has achieved the state-of-the-art performance on WebQSP and CWQ, confirming the effectiveness of diagnosing and remedying representation deficiencies. These results not only demonstrate the practical advantages of our method, but also highlight the fundamental importance of semantic representation robustness in complex reasoning tasks with LLMs.

540 Limitation

541 Although our proposed DR² method demonstrates
542 strong performance in the KGQA task by
543 mitigating the issue of non-deterministic reasoning,
544 we acknowledge several limitations in the
545 present work that point to directions for future
546 improvement:

547 **1) Dependency on Labeled Data:** While
548 our proposed DR² method effectively diagnoses
549 representation deficiencies, the process requires a
550 substantial amount of labeled data. Specifically,
551 identifying deficiencies requires sampling multiple
552 reasoning paths for each query and comparing
553 them with the annotated path. This reliance on
554 annotated data to locate semantic ambiguities limits
555 the method’s scalability in low-resource scenarios.
556 In future work, we will explore pathways to
557 detect representation deficiencies directly from all
558 relations in the knowledge graph, thereby reducing
559 the reliance on labeled data and enabling a more
560 automated diagnosis.

561 **2) Limited Validation on Task Generalization:**
562 Our current work evaluates the proposed diagnostic
563 and remediation method primarily within the
564 KGQA task, which validates its capability to
565 resolve semantic ambiguities in this setting.
566 However, its effectiveness on other tasks where
567 LLMs exhibit similar representation deficiencies
568 remains unexplored. For instance, in tasks
569 like Code Generation, Relation Extraction and
570 Sentiment Analysis, models may also struggle
571 with semantically similar options, leading to non-
572 deterministic reasoning. It remains an open
573 question whether the proposed deficiency diagnosis
574 and remediation method can transfer to these tasks.
575 In future work, we will explore the adaptability of
576 our approach across a wider range of tasks.

577 References

578 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda
579 Askeell, Anna Chen, Nova DasSarma, Dawn Drain,
580 Stanislav Fort, Deep Ganguli, Tom Henighan, et al.
581 2022. Training a helpful and harmless assistant with
582 reinforcement learning from human feedback. *arXiv*
583 *preprint arXiv:2204.05862*.

584 Ranran Bu, Jian Cao, Jianqi Gao, Shiyong Qian,
585 and Hongming Cai. 2025. [KaeDe: Progressive
586 generation of logical forms via knowledge-aware
587 question decomposition for improved KBQA](#). In
588 *Findings of the Association for Computational
589 Linguistics: EMNLP 2025*, pages 10958–10973,

Suzhou, China. Association for Computational
Linguistics. 590 591

Jianqi Gao, Jian Cao, Ranran Bu, Nengjun Zhu, Wei
Guan, and Hang Yu. 2025. Promoting knowledge
base question answering by directing llms to generate
task-relevant logical forms. In *Proceedings of
the AAAI Conference on Artificial Intelligence*,
volume 39, pages 23914–23922. 592 593 594 595 596 597

Yu Gu, Xiang Deng, and Yu Su. 2023. Don’t generate,
discriminate: A proposal for grounding language
models to real-world environments. In *Proceedings
of the 61st annual meeting of the association for
computational linguistics (volume 1: long papers)*,
pages 4928–4949. 598 599 600 601 602 603

Jiwoo Hong, Noah Lee, and James Thorne. 2024.
[ORPO: Monolithic preference optimization without
reference model](#). In *Proceedings of the 2024
Conference on Empirical Methods in Natural
Language Processing*, pages 11170–11189, Miami,
Florida, USA. Association for Computational
Linguistics. 604 605 606 607 608 609 610

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. 611 612 613 614 615 616 617 618

Jinhao Jiang, Kun Zhou, Wayne Xin Zhao, Yang Song,
Chen Zhu, Hengshu Zhu, and Ji-Rong Wen. 2025.
[KG-agent: An efficient autonomous agent framework
for complex reasoning over knowledge graph](#). In
*Proceedings of the 63rd Annual Meeting of the
Association for Computational Linguistics (Volume
1: Long Papers)*, pages 9505–9523, Vienna, Austria.
Association for Computational Linguistics. 619 620 621 622 623 624 625 626

Jinhao Jiang, Kun Zhou, Wayne Xin Zhao, and Ji-Rong
Wen. 2022. Unikgqa: Unified retrieval and reasoning
for solving multi-hop question answering over
knowledge graph. *arXiv preprint arXiv:2212.00959*. 627 628 629 630

Yunshi Lan and Jing Jiang. 2020. [Query graph gen-
eration for answering multi-hop complex questions
from knowledge bases](#). In *Proceedings of the 58th
Annual Meeting of the Association for Computational
Linguistics*, pages 969–974, Online. Association for
Computational Linguistics. 631 632 633 634 635 636

Tianle Li, Xueguang Ma, Alex Zhuang, Yu Gu, Yu Su,
and Wenhui 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. 637 638 639 640 641 642

Xiao Long, Liansheng Zhuang, Aodi Li, Minghong
Yao, and Shafei Wang. 2025. Eperm: An evidence
path enhanced reasoning model for knowledge graph
question and answering. In *Proceedings of the AAAI* 643 644 645 646

647		<i>Conference on Artificial Intelligence</i> , volume 39, pages 12282–12290.	
648			
649	Haoran Luo, E Haihong, Zichen Tang, Shiyao Peng,		
650	Yikai Guo, Wentai Zhang, Chenghao Ma, Guanting		
651	Dong, Meina Song, Wei Lin, et al. 2024. Chatkbqa:		
652	A generate-then-retrieve framework for knowledge		
653	base question answering with fine-tuned large		
654	language models. In <i>Findings of the Association for</i>		
655	<i>Computational Linguistics ACL 2024</i> , pages 2039–		
656	2056.		
657	LINHAO LUO, Yuan-Fang Li, Reza Haf, and Shirui		
658	Pan. 2024. Reasoning on graphs: Faithful and		
659	interpretable large language model reasoning. In		
660	<i>The Twelfth International Conference on Learning</i>		
661	<i>Representations</i> .		
662	Yu Meng, Mengzhou Xia, and Danqi Chen. 2024.		
663	SimpO: Simple preference optimization with		
664	a reference-free reward. <i>Advances in Neural</i>		
665	<i>Information Processing Systems</i> , 37:124198–124235.		
666	Alexander Miller, Adam Fisch, Jesse Dodge, Amir-		
667	Hossein Karimi, Antoine Bordes, and Jason Weston.		
668	2016. Key-value memory networks for directly		
669	reading documents . In <i>Proceedings of the 2016</i>		
670	<i>Conference on Empirical Methods in Natural</i>		
671	<i>Language Processing</i> , pages 1400–1409, Austin,		
672	Texas. Association for Computational Linguistics.		
673	Zhijie Nie, Richong Zhang, Zhongyuan Wang, and		
674	Xudong Liu. 2024. Code-style in-context learning		
675	for knowledge-based question answering. In		
676	<i>Proceedings of the AAAI Conference on Artificial</i>		
677	<i>Intelligence</i> , volume 38, pages 18833–18841.		
678	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,		
679	Carroll Wainwright, Pamela Mishkin, Chong Zhang,		
680	Sandhini Agarwal, Katarina Slama, Alex Ray,		
681	et al. 2022. Training language models to follow		
682	instructions with human feedback. <i>Advances in</i>		
683	<i>neural information processing systems</i> , 35:27730–		
684	27744.		
685	Rafael Rafailov, Archit Sharma, Eric Mitchell,		
686	Christopher D Manning, Stefano Ermon, and Chelsea		
687	Finn. 2023. Direct preference optimization: Your		
688	language model is secretly a reward model. <i>Advances</i>		
689	<i>in neural information processing systems</i> , 36:53728–		
690	53741.		
691	Yuan Sui, Yufei He, Nian Liu, Xiaoxin He, Kun Wang,		
692	and Bryan Hooi. 2025. FiDeLiS: Faithful reasoning		
693	in large language models for knowledge graph		
694	question answering. In <i>Findings of the Association</i>		
695	<i>for Computational Linguistics: ACL 2025</i> , pages		
696	8315–8330, Vienna, Austria. Association for		
697	Computational Linguistics.		
698	Haitian Sun, Bhuwan Dhingra, Manzil Zaheer, Kathryn		
699	Mazaitis, Ruslan Salakhutdinov, and William Cohen.		
700	2018. Open domain question answering using early		
701	fusion of knowledge bases and text. In <i>Proceedings</i>		
702	<i>of the 2018 conference on empirical methods in</i>		
703	<i>natural language processing</i> , pages 4231–4242.		
	Jiashuo Sun, Chengjin Xu, Lumingyuan Tang, Saizhuo		704
	Wang, Chen Lin, Yeyun Gong, Lionel Ni, Heung-		705
	Yeung Shum, and Jian Guo. 2024. Think-		706
	on-graph: Deep and responsible reasoning of		707
	large language model on knowledge graph. In		708
	<i>The Twelfth International Conference on Learning</i>		709
	<i>Representations</i> .		710
	Alon Talmor and Jonathan Berant. 2018. The		711
	web as a knowledge-base for answering complex		712
	questions. In <i>Proceedings of the 2018 Conference</i>		713
	<i>of the North American Chapter of the Association</i>		714
	<i>for Computational Linguistics: Human Language</i>		715
	<i>Technologies, Volume 1 (Long Papers)</i> , pages 641–		716
	651.		717
	Xiaqiang Tang, Jian Li, Nan Du, and Sihong Xie.		718
	2025. Adapting to non-stationary environments:		719
	Multi-armed bandit enhanced retrieval-augmented		720
	generation on knowledge graphs. In <i>Proceedings</i>		721
	<i>of the AAAI Conference on Artificial Intelligence</i> ,		722
	volume 39, pages 12658–12666.		723
	Hugo Touvron, Louis Martin, Kevin Stone, Peter		724
	Albert, Amjad Almahairi, Yasmine Babaei, Nikolay		725
	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti		726
	Bhosale, et al. 2023. Llama 2: Open foundation		727
	and fine-tuned chat models. <i>arXiv preprint</i>		728
	<i>arXiv:2307.09288</i> .		729
	Derong Xu, Xinhang Li, Ziheng Zhang, Zhenxi Lin,		730
	Zhihong Zhu, Zhi Zheng, Xian Wu, Xiangyu		731
	Zhao, Tong Xu, and Enhong Chen. 2025a.		732
	Harnessing large language models for knowledge		733
	graph question answering via adaptive multi-aspect		734
	retrieval-augmentation. In <i>Proceedings of the AAAI</i>		735
	<i>Conference on Artificial Intelligence</i> , volume 39,		736
	pages 25570–25578.		737
	Mufan Xu, Gewen Liang, Kehai Chen, Wei Wang, Xun		738
	Zhou, Muyun Yang, Tiejun Zhao, and Min Zhang.		739
	2025b. Memory-augmented query reconstruction for		740
	LLM-based knowledge graph reasoning . In <i>Findings</i>		741
	<i>of the Association for Computational Linguistics:</i>		742
	<i>ACL 2025</i> , pages 24068–24084, Vienna, Austria.		743
	Association for Computational Linguistics.		744
	Michihiro Yasunaga, Hongyu Ren, Antoine Bosselut,		745
	Percy Liang, and Jure Leskovec. 2021. QA-GNN:		746
	Reasoning with language models and knowledge		747
	graphs for question answering . In <i>Proceedings of</i>		748
	<i>the 2021 Conference of the North American Chapter</i>		749
	<i>of the Association for Computational Linguistics:</i>		750
	<i>Human Language Technologies</i> , pages 535–546,		751
	Online. Association for Computational Linguistics.		752
	Xi Ye, Semih Yavuz, Kazuma Hashimoto, Yingbo Zhou,		753
	and Caiming Xiong. 2022. RNG-KBQA: Generation		754
	augmented iterative ranking for knowledge base		755
	question answering . In <i>Proceedings of the 60th</i>		756
	<i>Annual Meeting of the Association for Computational</i>		757
	<i>Linguistics (Volume 1: Long Papers)</i> , pages 6032–		758
	6043, Dublin, Ireland. Association for Computational		759
	Linguistics.		760

761 Wen-tau Yih, Matthew Richardson, Christopher Meek,
 762 Ming-Wei Chang, and Jina Suh. 2016. The value
 763 of semantic parse labeling for knowledge base
 764 question answering. In *Proceedings of the 54th
 765 Annual Meeting of the Association for Computational
 766 Linguistics (Volume 2: Short Papers)*, pages 201–
 767 206.

768 Donghan Yu, Sheng Zhang, Patrick Ng, Henghui
 769 Zhu, Alexander Hanbo Li, Jun Wang, Yiqun Hu,
 770 William Yang Wang, Zhiguo Wang, and Bing Xiang.
 771 2022. Decaf: Joint decoding of answers and
 772 logical forms for question answering over knowledge
 773 bases. In *The Eleventh International Conference on
 774 Learning Representations*.

775 Jing Zhang, Xiaokang Zhang, Jifan Yu, Jian Tang, Jie
 776 Tang, Cuiping Li, and Hong Chen. 2022. Subgraph
 777 retrieval enhanced model for multi-hop knowledge
 778 base question answering. In *Proceedings of the 60th
 779 Annual Meeting of the Association for Computational
 780 Linguistics (Volume 1: Long Papers)*, pages 5773–
 781 5784.

782 Xinghua Zhang, Haiyang Yu, Cheng Fu, Fei Huang,
 783 and Yongbin Li. 2025. [IOPO: Empowering
 784 LLMs with complex instruction following via input-
 785 output preference optimization](#). In *Proceedings of
 786 the 63rd Annual Meeting of the Association for
 787 Computational Linguistics (Volume 1: Long Papers)*,
 788 pages 22185–22200, Vienna, Austria. Association
 789 for Computational Linguistics.

790 A Error Analysis of MemQ

791 In order to understand the limitations of the strong
 792 baseline method MemQ, we conducted a manual
 793 error analysis on its incorrect predictions within
 794 the first 500 samples of WebQSP dataset. The
 795 observed errors are systematically categorized
 796 into four primary types: **1) Relation Confusion**.
 797 It predicts a knowledge graph relation that is
 798 semantically related to the correct one; **2) Rare
 799 Relation**. The target relation is highly sparse
 800 or unseen in the training data; **3) Constraint
 801 Violation**. It incorrectly filters out valid answers
 802 by over-applying constraints; **4) Query Ambiguity**.
 803 The input natural language query is ambiguous or
 804 underspecified. As illustrated in Figure 7, Relation
 805 Confusion is the most prevalent failure, accounting
 806 for 53% of the analyzed errors. This suggests that
 807 the model’s primary weakness lies in distinguishing
 808 between semantically similar relations.

809 B Limitation of SFT

810 To empirically investigate the limitations of
 811 supervised fine-tuning (SFT) in discriminating
 812 closely related relations, we design a controlled
 813 experiment. We construct a dedicated dataset

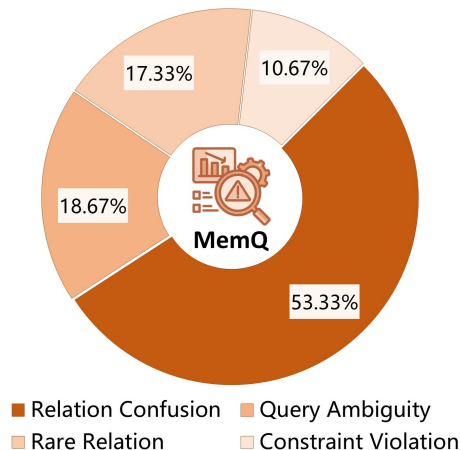


Figure 7: Error Analysis of MemQ.

814 where, for each query, the golden reasoning path
 815 must be exactly one of the two semantically
 816 similar relations: "language spoken" or "official
 817 language". This setup isolates the model’s ability
 818 to discriminate between easily confusable relations
 819 without interference from other error types.

820 We randomly split this dedicated dataset into
 821 training and test sets, with the results shown
 822 in Figure 8. The performance is measured by
 823 the average probability it assigns to generating
 824 the correct golden reasoning path. Although the
 825 SFT model achieves near-perfect accuracy on the
 826 training set, its performance drops to 75% on the
 827 test set. This discrepancy indicates that the SFT
 828 objective leads the model to primarily memorize
 829 surface patterns rather than learn to discriminate
 830 between closely related relations.

831 C Evaluation Metrics

832 In this section, we present the mathematical
 833 formulations and explanations for the primary
 834 metrics used in our evaluation. All reported results
 835 are averaged values.

836 **Hits@1**. Hits@1 quantifies the proportion of
 837 questions for which the top-ranked answer in the
 838 model’s output is correct. Let *Answer* represent
 839 the list of predicted answers, *Golden* denote the
 840 list of ground truth answers, and *total* represent
 841 the total number of questions in the dataset. The
 842 formula of Hits@1 is defined as follows:

$$843 \text{Hits@1} = \frac{\text{count}(\text{Answer}[0] \in \text{Golden})}{\text{total}}. \quad (8)$$

844 **Hits**. The Hits metric measures whether **any** of the
 845 model’s predicted answers is present in the set of

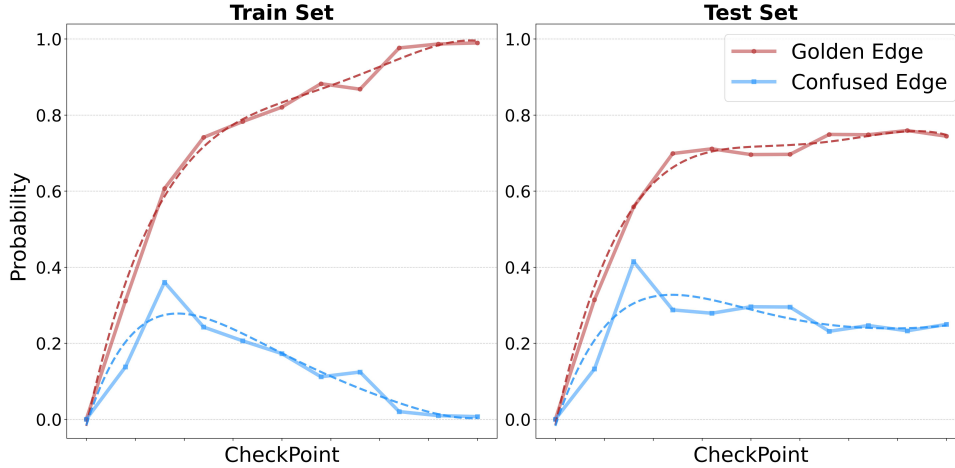


Figure 8: Limitation of SFT.

ground truth answers. It is formulated as:

$$Hits = \frac{\text{count}(\text{Answer} \cap \text{Golden} \neq \emptyset)}{\text{total}}. \quad (9)$$

Unlike Hits@1, which is a stricter measure requiring the first prediction to be correct, Hits considers the prediction successful if any of the returned answers matches a golden answer.

F1. The F1 score provides a balanced measure of the model’s overall answer quality by considering both precision and recall. We report the overall F1 score. This is calculated by first averaging the precision and recall values across all samples and then computing the F1 score based on them.

D Prompt Templates

The prompt templates used in our experiments are detailed below. We design two specific templates: one for the Semantic Distinguishability Evaluation (see Table 8) and another for the Inverse Question Generation via abductive reasoning (see Table 9).

E Further Representation Analysis

We provide additional evidence to illustrate how DR² mitigates representation deficiencies. This appendix presents extended analyses and visualizations on the identified defect set T_{Defect} , comparing the internal representations of DR² against those of the MemQ baseline. Figure 11 further demonstrate DR²’s effectiveness in separating previously confused relations.

Compared to MemQ, DR² consistently produces more distinct and better-separated representation clusters for the two reasoning structures within each confusable pair, with a clearer boundary

between them. This clear separation observed in the visualizations is supported by the quantitative results. For every defect pair examined, DR² achieves a significantly higher Separability Score ($SepSc$) than MemQ. This combined evidence robustly demonstrates that our method effectively enhances the model’s ability to discriminate between previously confusable reasoning structures.

F Extended Case Studies

To further illustrate how DR² enhances deterministic reasoning, we present the following extended case studies.

In Figure 9, the query is “*what country does rafael nadal play for?*”. The baseline SFT method (MemQ) exhibits confusion among multiple semantically related knowledge graph relations: “*nationality of the person someone*”, “*the country of which someone is a notable person*”, and “*the country of the country where someone has lived*”. exhibiting non-deterministic reasoning. In contrast, DR² maintains deterministic and correct reasoning reasoning, successfully selecting the precise relation “*nationality of the person someone*”. This case highlights DR²’s enhanced ability to resolve multi-relation confusion through improved semantic discrimination.

In Figure 10, the query is “*what region of the world is egypt associated with?*”. Although the baseline MemQ correctly predicts the relation “*the location that contains somewhere*”, it fails to apply the necessary constraint filtering for the “*notable type*”. In contrast, DR² accurately perceives the contextual cue “*region*” in the query and accordingly applies the necessary “*notable type*”

911 constraint during reasoning. This demonstrates
912 DR²'s enhanced ability to handle constraint-aware
913 reasoning, effectively resolving relation confusion
914 that arises within filtering constraints.

915 **G Details of Datasets**

916 In this section, we provide details about the datasets
917 utilized in this paper in Table 7. Both datasets are
918 publicly available and pose no security or privacy
919 concerns.

Dataset	Train	Test
WebQSP	3098	1639
CWQ	27639	3531

Table 7: Details about datasets.

920 **H Use of AI in Writing**

921 During the writing process, AI assistants were ex-
922 clusively used for text editing purposes, including
923 grammar correction and wording refinement.

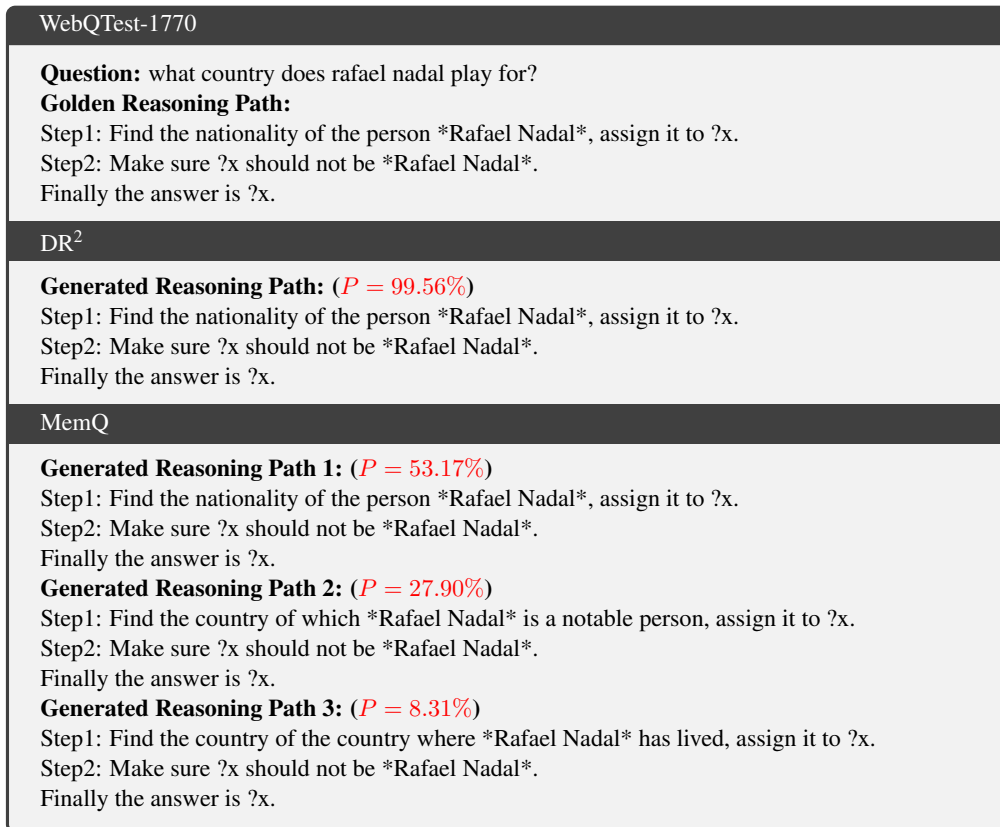


Figure 9: Compare a case of DR² and MemQ.

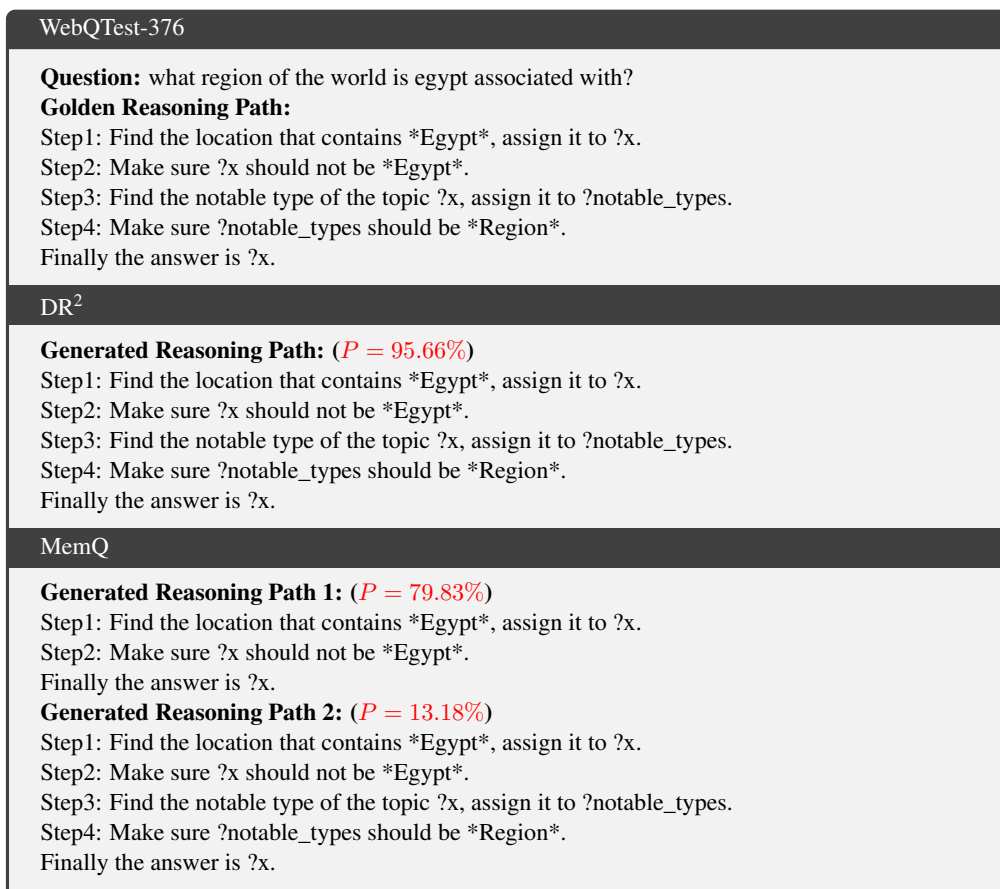


Figure 10: Compare a case of DR² and MemQ.

Semantic Distinguishability Evaluation

You are given two groups of questions:

Group1:

Group1 Hints: {group1_structure}

Group1 Questions:

{group1_question1}

{group1_question2}

{group1_question3}

{group1_question4}

Group2:

Group2 Hints: {group2_structure}

Group2 Questions:

{group2_question1}

{group2_question2}

{group2_question3}

{group2_question4}

Here is another question, your task is to classify it into either Group1 or Group2.

Respond strictly with either "Group1" or "Group2". Do not provide any explanations or other text.

Question:

question

Table 8: Semantic Distinguishability Evaluation

Inverse Question Generation

Based on the examples below, generate ONLY the natural language question based on the provided Entity and its Search Plan. The question should accurately reflect the search intent, mimicking the style and phrasing of the examples provided below.

Output NOTHING BUT the Question itself.

Here are the examples:

EXAMPLE1

Entity: {topic_entity1}

Search Plan: {golden_plan1}

Question: {question1}

EXAMPLE2

Entity: {topic_entity2}

Search Plan: {golden_plan2}

Question: {question2}

EXAMPLE3

Entity: {topic_entity3}

Search Plan: {golden_plan3}

Question: {question3}

Your Task:

Entity: {topic_entity}

Search Plan: {generated_plan}

Question:

Table 9: Inverse Question Generation

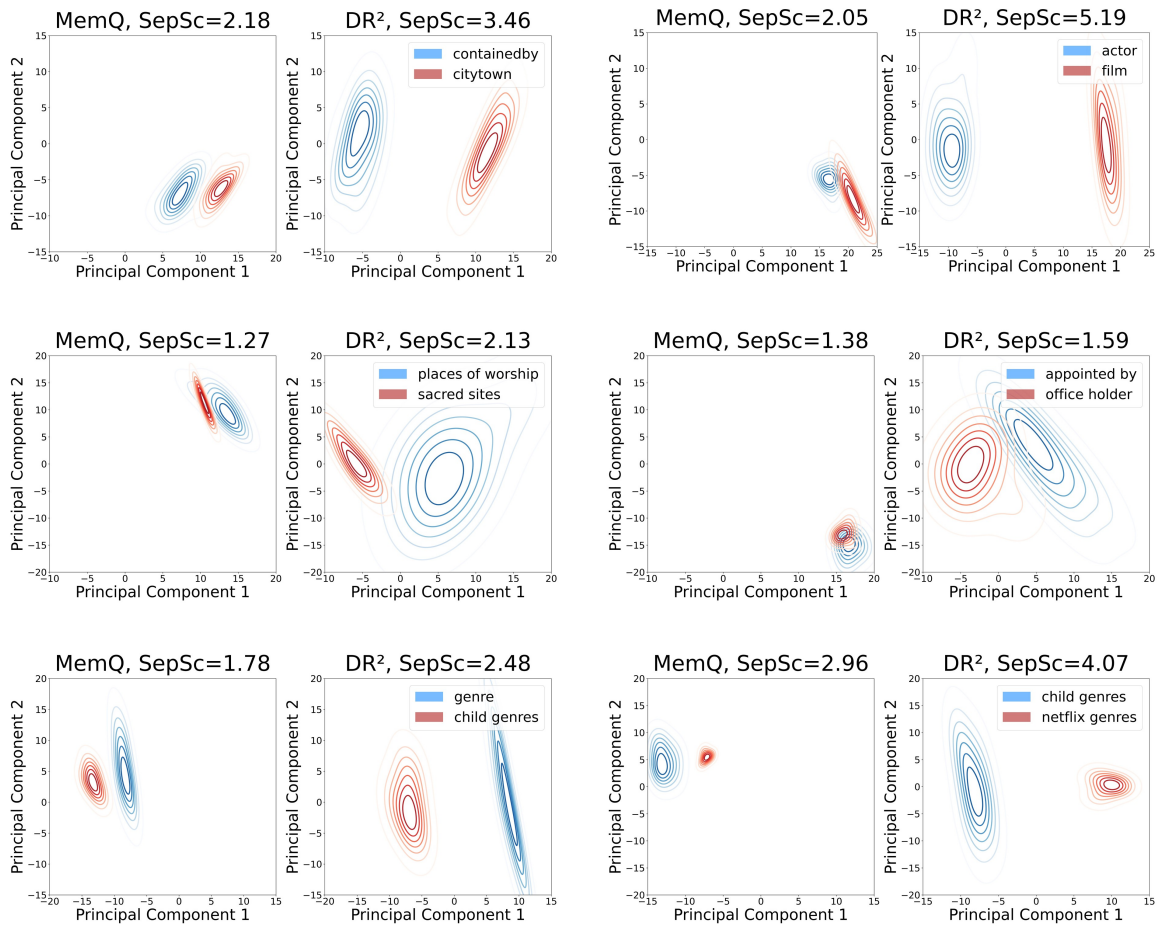


Figure 11: Further Representation Analysis.