Can KBQA Models Predict Their Reasoning Paths? Isomorphism Prediction Task as a Proxy

Anonymous ACL submission

Abstract

Despite achieving correct answers, we find that existing Knowledge Base Question Answering (KBQA) models struggle to follow the expected reasoning structures. We introduce the task of isomorphism prediction to enhance reasoning fidelity beyond answer generation, with a focus on generalization. We propose a contrastive knowledge co-distillation framework that unifies textual and graphical KBQA paradigms, improving overall isomorphism prediction and model generalization. Furthermore, incorporating isomorphism prediction as an auxiliary task could also improve KBQA performance.

006

017

018

1 Introduction and Related Work

The task of question answering over knowledge bases (KBQA) involves reasoning over structured sources of knowledge in the form of knowledge bases (KB) to answer natural language queries. Beyond improved answer accuracy, a key challenge in KBQA lies in understanding how the models perform and ensuring that they faithfully reconstruct the reasoning process. To that end, recent work has leveraged the idea of isomorphisms (Dutt et al., 2023) to characterize the complexity of KBQA questions. Isomorphisms act as a structural proxy for reasoning difficulty by grouping instances that exhibit similar reasoning patterns over the knowledge base. Prior work has explored using isomorphisms as a diagnostic test to investigate the generalization capabilities of KBQA systems. For example, Dutt et al. shows that leveraging gold isomorphisms as inference-time scaffolds improves zero-shot generalization without retraining.

In this work, we introduce the task of **isomorphism prediction** to improve reasoning fidelity in KBQA. Our task formulation is motivated by the observation, that when optimized for answer prediction, KBQA systems are able to generate spurious logical forms that do not conform to the underlying reasoning path but can lead to partial correct answers (Table 3). Furthermore, we observe that predicting the correct isomorphism category is challenging even for large language models (LLMs) (Table 6), highlighting the fact that the task requires models to capture structural dependencies beyond surface-level answer generation. Rather than solely using isomorphisms as a diagnostic tool, we frame them as a learning objective to encourage models to explicitly predict their underlying reasoning structures. 040

041

042

045

046

047

048

051

055

059

060

061

062

063

064

065

066

067

068

069

070

071

An advantage of this formulation is that it is applicable to both major KBQA paradigms: (i) semantic parsing-based approaches, which translate natural language queries into logical forms (e.g., S-expressions or SPARQL) for execution over the KB (Xie et al., 2022; Ye et al., 2021; Li et al., 2024), and (ii) information retrieval-based approaches, where models directly interact with the knowledge graph to retrieve answers (Das et al., 2022; Dutt et al., 2022; He et al., 2021). Building on this, we also propose a contrastive knowledge co-distillation framework that unifies these two paradigms to enhance isomorphism prediction.

Our experiments show that multitask learning with isomorphism prediction improves both KBQA and isomorphism prediction performance. Additionally, the proposed knowledge co-distillation framework bridges the strengths of both KBQA paradigms and enables better generalization.

2 Preliminaries

2.1 Isomorphism Prediction

We introduce isomorphism prediction to characterize reasoning paths following the definitions in Dutt et al. (2023). Each subgraph G_i represents the logical structure required to answer a question Q_i , where nodes correspond to entities and edges represent relations. Two subgraphs G_i and G_j are considered isomorphic if there exists a bijective

094

098

102

105

106

107

108

109

mapping $\psi: V_i \to V_j$ between their node sets and preserves structural adjacency:

$$(m,n) \in E_i \Leftrightarrow (\psi(m),\psi(n)) \in E_j$$

By assigning each subgraph to an isomorphism category C_i , we abstract away entity-specific details and focus purely on the structural reasoning pattern used to derive answers. We present the definitions and examples of each isomorphism type in Table 5.

Models are trained using a multi-class classification objective. We assess isomorphism prediction performance with macro F1-scores.

2.2 KBQA Tasks

Text Model: S-expression Generation Following Xie et al. (2022), the input to the text model consists of the question Q_i and the linearized representation of the subgraph (upper-left in Figure 1). The model generates an S-expression that retrieves the predicted answers when executed on the KB.

Graph Model: Node Classification The graph model operates directly on the subgraph G_i . It assigns probabilities to all nodes in the subgraph, indicating their likelihood of being answers. Training is optimized with binary cross-entropy.

KBQA Evaluation Mechanism We evaluate the aforementioned KBQA tasks with Hits@K, where K is the number of gold answers for a given question. This measures the proportion of correct answers in the top-K predictions.

For the text model, since S-expression generation does not produce ranked outputs, we use beam search to generate N_{beam} S-expression candidates, execute them through the KB, rank the executed answers by frequency, and compute Hits@K likewise. Refer to Section A.2 for detailed equations.

3 **Contrastive Knowledge Co-Distillation** for Isomorphism Prediction

Our Contrastive Knowledge Co-Distillation framework (Figure 1) consists of two key objectives: isomorphism prediction augmentation and contrastive representation alignment.

3.1 Isomorphism Prediction Augmentation

We employ two parallel encoding pathways. The 110 textual encoder produces a pooled embedding h_t 111 by processing the question along with a linearized 112 subgraph. The graph encoder, implemented as 113

a GNN, directly operates on the structured subgraph and generates a pooled graph-level representation h_a . These representations are concatenated as $h_{concat} = [h_t;, h_q]$ and passed through a classifier optimized via cross-entropy loss:

$$\mathcal{L}_{\rm iso} = -\sum_{i} \log P(C_i \mid h_{concat}) \qquad (1)$$

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

Contrastive Knowledge Co-Distillation 3.2

Unlike traditional one-way knowledge distillation, Contrastive Co-Distillation (Yao et al., 2024; Nourbakhsh et al., 2024) (CoD) fosters bidirectional knowledge transfer between text and graph models by contrastive representation learning and stop gradient operation.

As Tian et al. (2022) suggests, contrastive representation learning captures structural information from the teacher's representation space:

$$l_{\rm cl}(t,s) = -\log \frac{e^{sim(t,s)/\tau}}{\sum_{q} \mathbb{1}_{[q \neq t]} e^{sim(t,q)/\tau}} \quad (2)$$

, where t and s are teacher and student representations, q indicates other representations from the training data, sim(.,.) is cosine similarity, τ is temperature.

Based on this, we first define MLP projection heads to map text and graph representations into a shared space: $z_t = MLP_t(h_t)$ and $z_q = MLP_t(h_q)$, respectively. The CoD loss is computed as:

$$\mathcal{L}_{\text{CoD}} = \frac{1}{2} \sum_{i} [l_{\text{cl}}(z_i^{\text{text}}, \hat{z}_i^{\text{graph}}) + l_{\text{cl}}(\hat{z}_i^{\text{text}}, z_i^{\text{graph}})]$$
(3)

, where î is the stop gradient operator (Chen and He, 2021) to set the input variable to a constant.

Putting these together, our final objective jointly performs mutual distillation and model optimization end-to-end through a single loss:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{iso}} + \mathcal{L}_{\text{CoD}} \tag{4}$$

Experiments 4

Dataset 4.1

We employ the WebQuestionsSP (WebQSP) dataset (Yih et al., 2016), a popular benchmark in English for KBQA. Specifically, we use the dataset of Xie et al. (2022) where each question is accompanied with (i) a corresponding subgraph of the Freebase knowledge base (Bollacker et al., 2008) where the answer resides, and (ii) a corresponding logical form in the form of S-expressions or



Figure 1: Our contrastive knowledge co-distillation framework. The T5 encoder processes a linearized knowledge subgraph representation, while the GNN directly operates on the KG. Their representations are concatenated for isomorphism classification and projected into a shared space for contrastive co-distillation.

SPARQL-query. Such a design enables us to evaluate the performance of KBQA systems from either a semantic parsing or information retrieval paradigm. Additionally, to investigate different levels of KBQA generalization, we use the approach of Jiang and Usbeck (2022) to obtain a dev or test split with equal proportion of i.i.d., compositional, and zero-shot examples. We present the statistics of our dataset in Table 1.

Code	Desc.	i.i.d.	Comp	Z.S.	Total
T-0	—	50.3	0.0	49.7	54.5
T-1	— ——	37.3	44.3	18.4	23.5
T-2	• -•	17.1	47.1	35.7	5.2
T-3		83.3	6.7	10.0	2.2
T-4		12.8	81.5	5.6	14.5
ALL		40.8	24.9	34.3	100.0

Table 1: Distribution of isomorphisms over the generalization splits (i.i.d., compositional (Comp), zero-shot (Z.S.)) of WebQSP.

4.2 Models

Our text model is based on T5 (Raffel et al., 2023), while our graph model is built using Relational Graph Convolutional Network (RGCN) layers (Schlichtkrull et al., 2017). To incorporate question context, we first encode the question using T5 and concatenate the resulting embedding with each node before passing through the GNN.

4.3 Experiments

We establish T5 and GNN baselines, each trained separately for isomorphism prediction and their respective KBQA tasks defined in Section 2.2. We evaluate our approach under two settings: 1) CoD framework for isomorphism prediction; 2) Multitask KBQA using isomorphism prediction as an auxiliary task. We report in Section 5 the average performance over three seeds.

Moreover, we show that isomorphism prediction is challenging through two diagnostic experiments.

Firstly, we evaluate several widely-used LLMs on the isomorphism prediction task using few-shot prompting. As shown in Table 6, the models struggle to reliably predict isomorphisms.

Further, we analyze whether optimizing for Sexpression generation inherently preserves isomorphism structures. Although isomorphism categories can be deterministically derived from Sexpressions, models like T5 are trained to optimize answer accuracy rather than faithfully reconstructing reasoning paths. As a result, they may reach correct answers through spurious reasoning rather than the intended structural pattern. Indeed, isomorphism prediction performance drops by 16% overall when inferred post-hoc from generated Sexpressions compared with being explicitly learned in our T5 baseline (Table 3). This highlights the importance of directly modeling isomorphisms beyond relying on answer-driven supervision alone.

5 Results and Analysis

Isomorphism Prediction with CoD Overall, CoD outperforms both baselines (ALL in Table 2). We further stratify questions along generalization level and isomorphism category.

For generalization, GNN excels in i.i.d. cases but suffers in generalization, while T5 struggles the most in compositional settings. Although CoD performs lower in the i.i.d. setting, it significantly improves generalization, especially in compositional cases. This suggests that our contrastive knowledge co-distillation enables better adaptation rather than

213

214

165

166

168

169

170

171

172

173

174

175

176

177

156

157

158

160

162

163



Figure 2: Macro-F1 performance of the models/settings for different generalization levels in WebQSP.

Code	Desc.	T5	GNN	CoD
T-0	—	89.1	<u>88.8</u>	85.9
T-1	••	<u>60.8</u>	58.8	68.3
T-2	••	38.9	44.1	<u>43.5</u>
T-3		<u>45.7</u>	41.1	59.7
T-4		50.3	59.4	<u>51.8</u>
	ALL	57.0	<u>58.4</u>	61.9

Table 2: Macro-F1 of different settings (T5, GNN, and CoD) over isomorphism categories in WebQSP. ALL refers to the entire dataset. Best performance in bold, second-best underlined.

memorizing dataset biases.

215

216

217

218

219

222

224

Across isomorphism types, all models perform well in T-0 category (single-hop retrieval). As reasoning complexity increases, different model strengths become more evident. T5 performs well in linear chains (T-0, T-1) but struggles with more complex structures, while GNN is better with graph-structured constraints (T-2, T-4) but limited with sequential dependencies (T-3). Notably, CoD significantly improves on T-1 and T-3 and shows moderate gains in T-2 and T-4, which indicates that the unification brings together the complimentarities of the two models.

Multitask with KBQA and IsomorphismPredictionTable 4 shows that incorporating iso-

morphism prediction improves both KBQA and
isomorphism prediction tasks compared to singletask baselines. Our preliminary result shows that
isomorphism prediction provides additional structural supervision, which may help models better
capture reasoning patterns beyond answer retrieval.

Code	Desc.	T5 (Sexp)	T5 (Iso Pred)
T-0	—	82.0	89.1
T-1	— —	50.3	60.8
T-2	——	37.3	38.9
T-3		35.9	45.7
T-4		40.1	50.3
	ALL	40.9	57.0

Table 3: F1 performance of T5 on isomorphism prediction when inferred from generated S-expressions versus explicitly predicted as a supervised task.

Mode	l Task	KBQA	Iso Pred
	KBQA only	50.7	-
T5	Iso Pred only	-	59.0
	Multitask	52.2	61.7
	KBQA only	54.6	-
GNN	Iso Pred only	-	59.4
	Multitask	55.3	64.0

Table 4: Comparison of the respective task baselines and the multitask setting using isomorphism prediction as an auxiliary task.

6 Conclusion

We introduce isomorphism prediction task to enhance reasoning fidelity in KBQA. Our contrastive knowledge co-distillation framework improves isomorphism prediction and generalization, particularly in compositional and zero-shot settings. Additionally, isomorphism prediction as an auxiliary task improves KBQA performance, suggesting structural reasoning signals could aid answer generation. Future work can explore broader model architectures and datasets.

236

237

238

239

240

241

242

243

244

245

247

248

249

250

251

252

253

254

255

256

7 Limitations

Model Scope We focus on T5 and GNN-based models; future work could extend to larger LLMs and alternative graph reasoning frameworks.

Dataset Diversity Our experiments use WebQSP. Future work could extend evaluations to benchmarks with more diverse KG schemas.

Explicit Isomorphism Learning Future work could explore unsupervised learning to infer reasoning structures without predefined labels.

References

- Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. 2008. Freebase: a collaboratively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data*, SIGMOD '08, page 1247–1250, New York, NY, USA. Association for Computing Machinery.
- Xinlei Chen and Kaiming He. 2021. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15750–15758.
- Rajarshi Das, Ameya Godbole, Ankita Naik, Elliot Tower, Manzil Zaheer, Hannaneh Hajishirzi, Robin Jia, and Andrew McCallum. 2022. Knowledge base question answering by case-based reasoning over subgraphs. In *International conference on machine learning*, pages 4777–4793. PMLR.
- Ritam Dutt, Kasturi Bhattacharjee, Rashmi Gangadharaiah, Dan Roth, and Carolyn Rose. 2022. Perkgqa: Question answering over personalized knowledge graphs. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 253–268.
- Ritam Dutt, Sopan Khosla, Vinayshekhar Bannihatti Kumar, and Rashmi Gangadharaiah. 2023. Grailqa++: A challenging zero-shot benchmark for knowledge base question answering. In *Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 897–909.
- Ritam Dutt, Dongfang Ling, Yu Gu, and Carolyn Penstein Rosé. Leveraging isomorphisms to facilitate zero-shot kbqa generalization.
- 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.
- Longquan Jiang and Ricardo Usbeck. 2022. Knowledge graph question answering datasets and their generalizability: Are they enough for future research? *arXiv* preprint arXiv:2205.06573.
- Chunhui Li, Yifan Wang, Zhen Wu, Zhen Yu, Fei Zhao, Shujian Huang, and Xinyu Dai. 2024. Multisql: A schema-integrated context-dependent text2sql dataset with diverse sql operations. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 13857–13867.
- Armineh Nourbakhsh, Zhao Jin, Siddharth Parekh, Sameena Shah, and Carolyn Rose. 2024. AliGATr: Graph-based layout generation for form understanding. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 13309– 13328, Miami, Florida, USA. Association for Computational Linguistics.

- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2023. Exploring the limits of transfer learning with a unified text-to-text transformer. *Preprint*, arXiv:1910.10683.
- Michael Schlichtkrull, Thomas N. Kipf, Peter Bloem, Rianne van den Berg, Ivan Titov, and Max Welling. 2017. Modeling relational data with graph convolutional networks. *Preprint*, arXiv:1703.06103.
- Yonglong Tian, Dilip Krishnan, and Phillip Isola. 2022. Contrastive representation distillation. *Preprint*, arXiv:1910.10699.
- Tianbao Xie, Chen Henry Wu, Peng Shi, Ruiqi Zhong, Torsten Scholak, Michihiro Yasunaga, Chien-Sheng Wu, Ming Zhong, Pengcheng Yin, Sida I. Wang, Victor Zhong, Bailin Wang, Chengzu Li, Connor Boyle, Ansong Ni, Ziyu Yao, Dragomir Radev, Caiming Xiong, Lingpeng Kong, Rui Zhang, Noah A. Smith, Luke Zettlemoyer, and Tao Yu. 2022. UnifiedSKG: Unifying and multi-tasking structured knowledge grounding with text-to-text language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 602–631, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Hao-Ren Yao, Luke Breitfeller, Aakanksha Naik, Chunxiao Zhou, and Carolyn Rose. 2024. Distilling multi-scale knowledge for event temporal relation extraction. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*, CIKM '24, page 2971–2980, New York, NY, USA. Association for Computing Machinery.
- Xi Ye, Semih Yavuz, Kazuma Hashimoto, Yingbo Zhou, and Caiming Xiong. 2021. Rng-kbqa: Generation augmented iterative ranking for knowledge base question answering. *arXiv preprint arXiv:2109.08678*.
- 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* (Volume 2: Short Papers), pages 201–206.

A Appendix

A.1 Isomorphism Examples

See Table 5.

A.2 Hits@K Computations

Let A_i be the gold answer set for a given question. The graph model ranks all nodes in the subgraph by predicted probabilities. Given a ranked list \hat{A}_i , for some $\varepsilon > 0$, Hits@ K_{graph} is computed as:

$$Hits@K_{graph} = \frac{|TopK(\hat{A}_i) \cap A_i|}{K + \varepsilon}$$
(5)

360

361

362

364

257

258

259

260

261

263

264

265

267

268

269

270

271

272

277

278

279

281

282

284

290

296

297

298

303

306

307

310

311

312

Iso-Type	Illustration	Definition	Example Question	S-expression
T-0	••	Direct 1-hop connection from constraint to answer	What is the name of money in Brazil?	(JOIN (R loca- tion.country.currency_used) m.015fr)
T-1		2-hop linear path	Where does the Queen of Den- mark live?	(JOIN (R people.place_lived.location) (JOIN (R people.person.places_lived) m.0g2kv))
T-2	• -•-•	V-pattern with two constraints meeting at a shared node	What was Elie Wiesel's fa- ther's name?	(AND (JOIN people.person.gender m.05zppz) (JOIN (R peo- ple.person.parents) m.02vsp))
T-3		A chain pattern connecting constraints serially	Where did Joe Namath attend college?	(AND (JOIN com- mon.topic.notable_types m.01y2hnl) (JOIN (R educa- tion.education.institution) (JOIN (R people.person.education) m.01p_3k)))
T-4		Y-pattern with merging con- straints	Who does Zach Galifianakis play in The Hangover?	(JOIN (R film.performance.character) (AND (JOIN film.performance.film m.0n3xxpd) (JOIN (R film.actor.film) m.02_0d2)))

Table 5: Isomorphism types with their corresponding definitions, example questions, and S-expressions.

Unlike the graph model, the text model does not inherently rank its predictions. To approximate a ranking mechanism, we employ beam search to generate N_{beam} candidate S-expressions $S_{i,j}, j =$ $1, ..., N_{\text{beam}}$. We then execute these S-expressions through KB to obtain a predicted answer set P_i , and aggregate $P_{i,j}$ by their frequency across all beams. Using this ranked set, for some $\varepsilon > 0$, Hits@ K_{text} is computed as:

367

369

371

373

374

376

377

378

396

$$Hits@K_{\text{text}} = \frac{|TopK(Rank(\cup^{N_{\text{beam}}}P_i)) \cap A_i|}{K + \varepsilon}$$
(6)

where $Rank(\cup^{N_{beam}}P_i)$ refers to the aggregated ranking of answer candidates obtained from executing S_i through KB.

A.3 Few-shot LLM on Isomorphism Prediction

We evaluate a couple of widely-used LLMs on the isomorphism prediction task with few-shot prompting, including GPT-3.5-turbo and GPT-4o-mini. As shown in Table 6, these models struggle to reliable predict isomorphisms even with multiple examples per type of isomorphism. We try k-shot prompting with k = 1, 3, where we include k examples of each isomorphism type (T-0 to T-4), selected randomly from the training split. The exact prompts used can be found in Appendix A.4.

We experiment with not only the number of fewshot examples provided to the model, but also the technique used to serialize the knowledge graph tuples into a text format as well as the level of detail in the prompt about descriptions of particular isomorphisms. For serializing the knowledge graph tuples of the form $(entity_1, rel, entity_2)$ we try:

1. **Basic serialization:** where we simply concatenate the knowledge graph tuples using whitespace, for example "entity₁ rel entity₂". 397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

2. **Descriptive serialization:** where we concatenate each individual tuple with slightly more description, for example "entity₁ is connected to entity₂ via relation rel".

We try two levels of isomorphism description detail in our prompt. In the first setting, Prompt 1 (Appendix A.4.1), we provide a brief textual description of each of the isomorphisms' structural characteristics. Whereas in Prompt 2 (Appendix A.4.2), we do not provide any description whatsoever of individual isomorphism categories. The LLMs' final answers are extracted using a regex expression to match the last occurrence of the pattern "T-X", which indicates the model's isomorphism prediction. These predictions are then evaluated using a standard macro F-1 score. These scores, across all experiments, are shown in Table 6. We find that even with few-shot examples, and across all our prompting methods described above, the best performance achieved is a mean macro F-1 of 0.15 by the gpt-3.5-turbo model when given 3 examples per isomorphism class, basic serialized tuples and brief descriptions of isomorphisms' structure.

A.4 Few-shot LLM Prompt

Below are the two versions of prompts we experimented with. Prompt 1 contains brief structural descriptions of each isomorphism category, whereas

Model Configuration	Macro F1
GPT-3.5-turbo	
k=1 (base)	0.09220
k=1 (descriptive tuples)	0.14062
k=1 (descriptive tuples, no iso. desc. in prompt)	0.14118
k=3 (base)	0.15474
GPT-4o-mini	
k=1 (base)	0.10530
k=3 (base)	0.10273
k=3 (descriptive tuples)	0.14423
k=3 (descriptive tuples, no iso. desc. in prompt)	0.11832
k=5 (base)	0.07076
k=5 (descriptive tuples)	0.12022
k=5 (descriptive tuples, no iso. desc. in prompt)	0.09122

Table 6: Isomorphism prediction performance of GPT-3.5-turbo and GPT-4o-mini using few-shot prompting. The base configuration refers to when we serialize in a basic manner and provide brief isomorphism descriptions in the prompt.

Prompt 2 simply instructs the model to identify the isomorphism based on the examples provided.

A.4.1 Prompt 1: Structural Descriptions

System prompt: "You are a helpful assistant that identifies isomorphism patterns in knowledge graphs."

User prompt: "Given a question, its entities, and knowledge graph tuples, determine the isomorphism pattern that shows how constraints connect to reach the answer node. In this classification: T-0 means a direct 1-hop connection from constraint to answer, T-1 is a 2-hop linear path, T-2 is a V-pattern with two constraints meeting at a shared node, T-3 is a chain pattern connecting constraints serially, T-4 is a Y-pattern with merging constraints, and T-5+ involve more complex multi-hop patterns. Analyze the structure by tracing the paths from constraints to the answer, counting hops and noting how paths merge or branch. Respond with "Isomorphism: T-X" where X is the pattern number (0-4), output only the final answer. Find some examples below: { Question: ...

Entities: (example serialized entities)

Serialized tuples from knowledge graph: (example serialized knowledge graph tuples) Isomorphism: T-X } (k examples for each type of

isomorphism)

###

Question: (target question) Entities: (target serialized entities)

Serialized tuples from knowledge graph: (target 460

serialized knowledge graph tuples) 461 Isomorphism:" 462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

491

492

493

494

495

A.4.2 Prompt 2: Non-descriptive instructions

System prompt: "You are a helpful assistant that identifies isomorphism patterns in knowledge graphs."

User prompt: "Given a question, its entities, and knowledge graph tuples, determine the isomorphism pattern that shows how constraints connect to reach the answer node. Analyze the structure by tracing the paths from constraints to the answer, counting hops and noting how paths merge or branch. Respond with "Isomorphism: T-X" where X is the pattern number (0-4), output only the final answer. Find some examples below: { Question: ...

Entities: (example serialized entities) Serialized tuples from knowledge graph: (example serialized knowledge graph tuples) Isomorphism: T-X } (k examples for each type of isomorphism)

###

GPU hours.

	-105
Question: (target question)	484
Entities: (target serialized entities)	485
Serialized tuples from knowledge graph: (target	486
serialized knowledge graph tuples)	
Isomorphism:"	488
A.5 Hyperparameter Settings	489
On average, our total experiments take around 15	490

A.5.1 **Experiments on Isomorphism Prediction with CoD**

We use the following hyperparameters to obtain results in Table 2.

Model	Batch	Dropout	Others	
T5	8	0.2	-	
GNN	10	0.2	-	
CoD	6	0.3	Weight Decay: 1e-3	
Shared Space Dim: 2048				
Input Max Length: 512, Patience: 5, LR: 5e-5				

Table 7: Hyperparameters used for experiments in Table 2. Results are averaged over three seeds.

448

449

450

451

452

453

454

455 456

457

458

459

429

430

A.5.2 Experiments on Multitask with KBQA and Isomorphism Prediction

496 497 498

499

We use the following hyperparameters to obtain results in Table 4.

Model	Batch	Dropout	Others
T5 Baseline	10	-	Generation Max Len: 128
GNN Baseline	6	0.2	-
T5 Multitask	10	0.3	Weight Decay: 1e-3 Generation Max Len: 128
GNN Multitask	4	0.2	-
Input Max Length: 512, Patience: 5, LR: 5e-5			

Table 8: Hyperparameters used for KBQA and multitask experiments in Table 4. Results are averaged over three seeds.

A.6 System Specifications

See Table 9.

Component	Specification
GPU	NVIDIA A100 80GB PCIe
CPU	AMD EPYC 7763 (256 vCPUs)
RAM	1TB
CUDA Version	12.6
GPU Memory	80GB

Table 9: Hardware specifications of the computationalresources used for experiments.

A.7 Potential Risks and Considerations

Our work builds on WebQSP and Freebase, which may inherit biases from their original data collection. While our focus is on structural reasoning rather than entity-specific biases, these biases could still affect model behavior. Additionally, although we do not train large models from scratch, prompting LLMs, fine-tuning T-5, and training GNN still lead to computational costs, contributing to the environmental footprint.

501

502

503

504

505

506

507

508

510