

Neuro-Symbolic Integration Brings Causal and Reliable Reasoning Proofs

Anonymous ACL submission

Abstract

Though prompting LLMs with various reasoning structures produces intermediate reasoning steps along with answers, these steps are not ensured to be causal and reliable due to the inherent defects of LLMs. Tracking such deficiencies, we present a neuro-symbolic integration framework, in which a neural LLM is used to represent the knowledge of the problem while an LLM-free symbolic solver is adopted to do deliberate reasoning using the knowledge. Specifically, customized meta-interpreters are implemented to generate intermediate reasoning proofs and to support various search strategies. These reasoning proofs are ensured to be causal and reliable because of the deterministic executing nature of the symbolic solvers. We conduct experiments on two logical reasoning and one arithmetic reasoning datasets. On ProofWriter, our method surpasses the CoT baseline by nearly double in reasoning accuracy and more than triple in reasoning proof similarity. On GSM8K, our method also shows accuracy improvements and nearly doubled proof similarity. Our code is released at <https://anonymous.4open.science/r/CaRing-477B>.

1 Introduction

Large language models (LLMs), like LLaMA-2 (Touvron et al., 2023) and GPT-4 (OpenAI, 2023), are shown to be effective on several reasoning tasks but still struggle with structurally complex reasoning problems, such as logical reasoning (Tafjord et al., 2021) and arithmetic reasoning (Cobbe et al., 2021; Ribeiro et al., 2023). To tap into the potential of LLMs for better complex reasoning, existing works primarily focus on iteratively prompting LLMs to search over reasoning structures such as chains (e.g., CoT) (Wei et al., 2022; Wang et al., 2023; Zhou et al., 2023), trees (e.g., Tree-of-Thoughts, RAP) (Yao et al., 2023; Long, 2023; Hao et al., 2023), and graphs (e.g.,

Graph-of-Thoughts) (Besta et al., 2023; Zhang et al., 2023; Sun et al., 2023).

Despite the effectiveness of such methods over various complex reasoning problems, it is observed that they often give correct results with erroneous intermediate steps (Ye and Durrett, 2022; Saparov and He, 2023; Ribeiro et al., 2023). For example, Ribeiro et al. (2023) showed that even though prompting GPT-3 given structured intermediate steps as demonstrations yields an average accuracy of 33.84% on five complex reasoning datasets, the average similarity between the predicted and the gold reasoning proofs is merely 0.72%. This discrepancy between reasoning accuracy and reasoning proof similarity raises pressing concerns about the reliability and causality of the underlying reasoning process in LLMs, as shown in Figure 1.

The discrepancies identified in the reasoning capabilities of LLMs underscore their limitations in emulating human-like deliberate reasoning. One natural solution could be adopting an LLM-free deliberate reasoning engine. Inspired by the seminal work of Kowalski (1979), which argued that a problem-solving algorithm benefits from separating the *Logic* component (i.e., the knowledge which can be used to solve the problem) and the *Control* component (i.e., the problem-solving strategy with which the knowledge can be used), we propose a neuro-symbolic integration approach consisting of two components: (1) LLM-based symbolic representation generator (SYMGEN; §3.1), which translates natural languages into formal knowledge representations that can be used for symbolic inference; (2) LLM-free symbolic inference engine (SYMINFER; §3.2), which performs deliberate reasoning by executing the symbolic representations. The execution strategy is implemented with our customized meta-interpreters, allowing (i) tracing of the reasoning process (§3.2.1); (ii) adoption of various search strategies (§3.2.2). Most importantly, by putting LLMs under quarantine during deliberate

reasoning, our approach produces reasoning traces that are strictly causal and immune from hallucinations.

To demonstrate its effectiveness in producing better reasoning proofs, we evaluate our CARING (Causal and Reliable Reasoning) framework on three reasoning datasets that contain reasoning proof annotations, including two logical reasoning datasets, ProofWriter (Tafjord et al., 2021) and PrOntoQA (Saparov and He, 2023), and one arithmetic reasoning dataset, GSM8K (Ribeiro et al., 2023). CARING consistently outperforms the CoT baseline and existing methods in terms of answer accuracy, reasoning proof similarity, and reasoning proof accuracy. On the challenging ProofWriter dataset, CARING using Code-LLaMA-34B yields an answer accuracy of 96.5% and a reasoning proof similarity of 81.0%, while previous SoTA achieved an answer accuracy of 79.7%. Further analysis indicates CARING remains robust when the reasoning problem becomes more complex.

Overall, our contributions in this paper include:

- As far as we know, CARING is the first LLM-based neuro-symbolic integration approach that customizes symbolic interpreters to generate reasoning proofs.
- We present an implementation using Prolog representations and conduct experiments on three datasets. Empirically, our framework gains significant improvements over strong baselines and existing methods with both final answers and reasoning proofs.

2 Related Work

2.1 Explainable Complex Reasoning

The Chain-of-Thought prompting method, which found out reasoning with LLMs benefits from generating intermediate steps, has sparked a recent trend in how to better do reasoning while remaining explainable. Several works investigated using other reasoning structures, such as trees (Yao et al., 2023; Long, 2023) and graphs (Besta et al., 2023; Zhang et al., 2023). These approaches have shown improved performance, particularly in complex reasoning tasks where the processes involved are often more intricate than simple linear chains. However, despite the alignment of their reasoning proof structures with the gold-standard proofs, these methods still face challenges in ensuring causality and reliability. This limitation stems from their reliance on

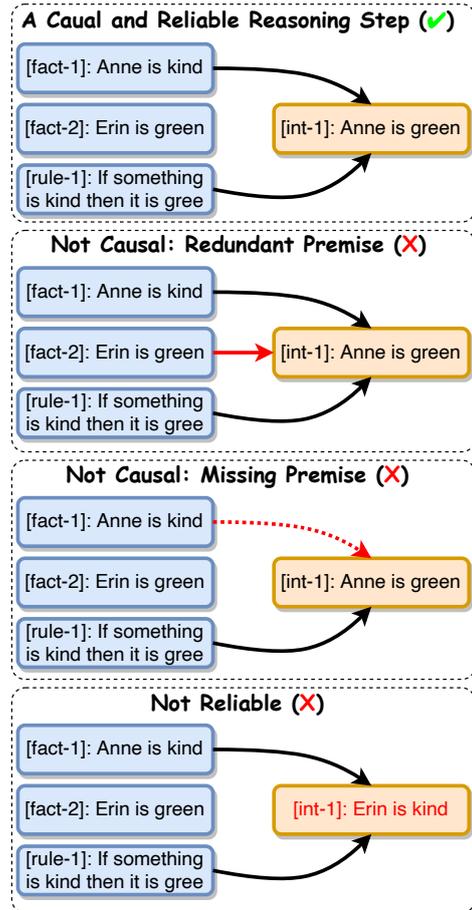


Figure 1: Illustrations of how causality and reliability play important roles in reasoning. LLMs may be (i) non-causal by selecting redundant premises or ignoring relevant ones and (ii) non-reliable by hallucinating incorrect contents during inference.

LLMs for deliberate reasoning, which are prone to hallucinations and may compromise causality.

Some other recent works adopted a less structured manner (Tafjord et al., 2022; Creswell et al., 2023; Kazemi et al., 2023). For example, Selection-Inference (Creswell et al., 2023) divides the reasoning process into two phases: (1) the Selection phase for selecting the premises that might be relevant for the next round of inference, and (2) the Inference phase for conducting a single reasoning step with the selected knowledge fragments.

2.2 Neuro-symbolic Reasoning

Neuro-symbolic systems attempt to leverage the strengths of both neural networks and symbolic reasoning (Andreas et al., 2016; Neelakantan et al., 2017; Hudson and Manning, 2019; Gupta et al., 2020; Nye et al., 2021). This includes the use of neural networks for pattern recognition and learning from unstructured data, integrated with sym-

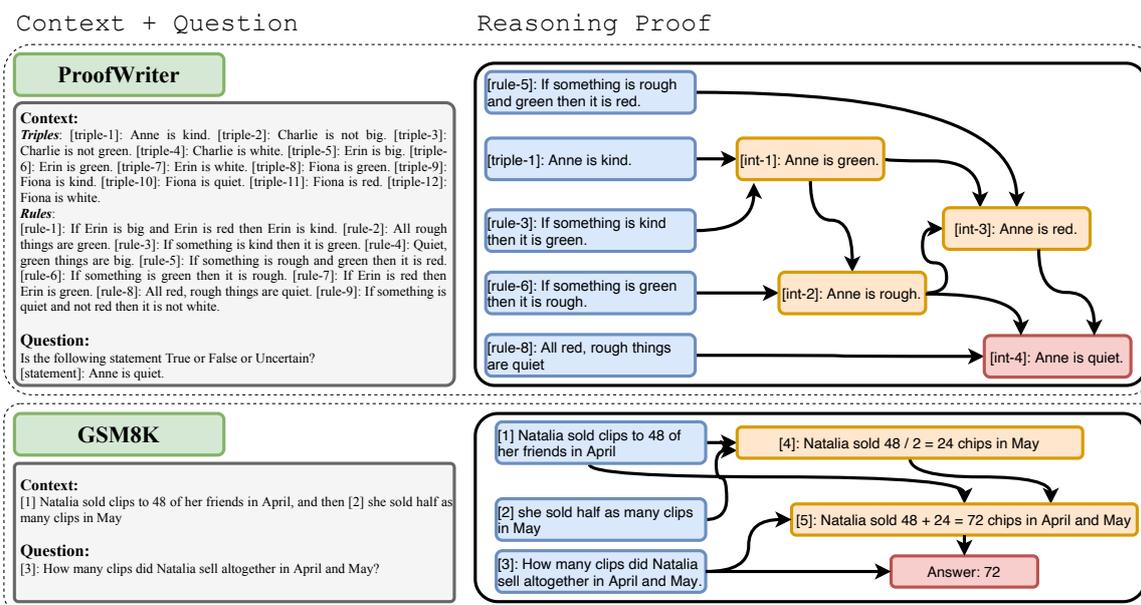


Figure 2: Two examples of complex/structured reasoning problems from ProofWriter and GSM8K, respectively. The reasoning proofs in such problems formulate directed acyclic graphs (DAGs) in a multi-step and multi-premise manner.

151 bolic systems for rule-based reasoning and knowl- 180
 152 edge representation. Despite significant progress, 181
 153 neuro-symbolic reasoning faces challenges, nota- 182
 154 bly in scalability and the efficient integration of 183
 155 learning and reasoning components. 184

156 Recent advancements in neuro-symbolic re- 185
 157 search, particularly in reasoning over text, have 186
 158 utilized LLMs to encapsulate knowledge from un- 187
 159 structured human languages, as noted in Lyu et al. 188
 160 (2023); Pan et al. (2023). These methods typically 189
 161 translate natural language into symbolic represen- 190
 162 tations for subsequent execution-based reasoning. 191
 163 However, they have not fully explored the capabili- 192
 164 ties of symbolic solvers in generating detailed rea- 193
 165 soning proofs. In contrast, our approach leverages 194
 166 customized meta-interpreters in conjunction with 195
 167 symbolic solvers to uncover and articulate the un- 196
 168 derlying reasoning proofs. This not only enhances 197
 169 the transparency of automatic reasoning systems 198
 170 but also simplifies the process for humans to verify 199
 171 their correctness and safety. 200

172 3 CARING

173 The problems we focus on are featured with struc- 201
 174 tured or complex reasoning. As depicted in Fig- 202
 175 ure 2, these problems typically necessitate multi- 203
 176 step and multi-premise reasoning over a directed 204
 177 acyclic graph (DAG), where individual nodes sig- 205
 178 nify distinct knowledge fragments and directed 206
 179 edges denote reasoning steps. Each reasoning 207
 208

180 step uses existing knowledge to infer new relevant 180
 181 knowledge. Numerous knowledge fragments are 181
 182 often aggregated to infer a new one, which we 182
 183 denote as “multi-premise”. The solver usually per- 183
 184 forms multiple such steps to reach an ultimate goal, 184
 185 which we denote as “multi-step”. This entire rea- 185
 186 soning process naturally composes a DAG. 186

187 We are interested in providing accurate answers 187
 188 along with causal and reliable explanations for such 188
 189 reasoning problems. This motivates our investiga- 189
 190 tion of LLM-free deliberate reasoning engines. 190
 191 The seminal work of Kowalski (1979) proposed 191
 192 that $Algorithm = Logic + Control$, where *logic* 192
 193 refers to the knowledge which can be used to solve 193
 194 the problem and *control* refers to the problem- 194
 195 solving strategy in which the knowledge can be 195
 196 used. They further proved that an algorithm ben- 196
 197 efits from separating the *logic* component and the 197
 198 *control* component. Inspired by this, we present 198
 199 CARING (Causal and Reliable Reasoning), a mod- 199
 200 ular approach consisting of two components: 200

- 201 • SYMGEN: LLM-based symbolic representa- 201
 202 tion generator (§3.1), which translates natural 202
 203 languages into formal symbolic knowledge 203
 204 representations that can be used for symbolic 204
 205 inference. A major difference between previ- 205
 206 ous work and our method is that we only use 206
 207 LLMs to represent knowledge but not to do 207
 208 deliberate reasoning. 208

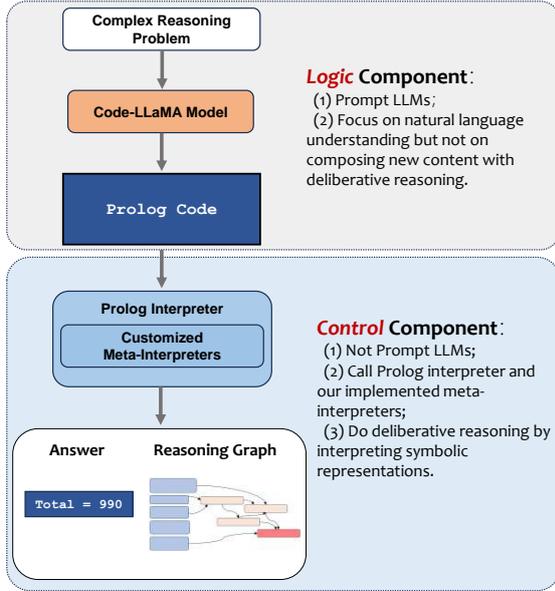


Figure 3: Illustration of our CARING framework, consisting of a *Logic* component and a *Control* component.

- **SYMINFER:** LLM-free symbolic inference engine (§3.2), which performs deliberate reasoning by executing the symbolic representations provided by SYMGEN. By implementing customized meta-interpreters, SYMINFER supports (i) causal and reliable tracing of the reasoning process (§3.2.1); (ii) various search strategies, such as Depth-First Search (DFS) and Iterative Deepening Search (IDS) (§3.2.2).

The execution-based tracing approach of SYMINFER guarantees both causality and reliability. Under the principle of **Causality**, the inference of a new knowledge piece is strictly linked to those existing fragments that are relevant, ensuring precise and limited attribution. This implies that a causal relationship is established only when the preceding event (at the base of the edge) directly influences the subsequent event (at the apex). Regarding **Reliability**, the content within each newly inferred node is the result of a deterministic process, safeguarding it from the kinds of erroneous hallucinations often encountered in outputs from LLMs.

3.1 SYMGEN: Symbolic Representation Generator

To represent *logic* (i.e., the knowledge which can be used to solve the problem), we adopt a popular logic programming language, Prolog (Colmerauer and Roussel, 1996). Prolog is a declarative pro-

Natural Language	Prolog Code
<i>Fiona is green.</i>	1 green(fiona).
<i>All red, rough things are quiet.</i>	1 quiet(X) :- 2 red(X), rough(X).
<i>Tina makes \$18.00 an hour.</i>	1 wage(18.00).
<i>(she is eligible for overtime,) which is paid by your hourly wage + 1/2 your hourly wage.</i>	1 overtime_wage(W) :- 2 wage(W1), 3 W is 1.5 * W1.

Table 1: Examples of natural languages and their Prolog representations. It can be seen that the Prolog code is highly declarative, so the LLM in SYMGEN is only required to do straightforward natural language understanding and translation but not reasoning. In other words, the LLM does not need to infer new knowledge, thus avoiding hallucination as much as possible.

gramming language, in which *logic* is expressed as relations (called Facts and Rules), with several examples shown in Table 1. A computation is initiated by running a query over these relations. We will delve into the computation of Prolog in §3.2.

Though LLMs are prone to hallucinate erroneous facts when composing new knowledge, they are shown to be powerful at understanding natural languages and directly translating them into other formats (Ye and Durrett, 2022; Saparov and He, 2023). To utilize such a strong point while avoiding the defect, we only use LLMs to translate natural languages into Prolog representations but not to do deliberate reasoning. Specifically, we few-shot prompt LLMs with several human-written in-context demonstrations, each containing a problem and corresponding Prolog representations, which are later used for symbolic inference.

3.2 SYMINFER: Symbolic Inference Engine

We use SYMINFER to produce answers and reasoning traces by executing the aforementioned symbolic representations. Since we adopt Prolog to represent knowledge, our symbolic inference engine is naturally instantiated with Prolog interpreters. By default, the SWI-Prolog (Wielemaker et al., 2012) interpreter adopts the Depth-First Search (DFS) backtracking strategy and does not yield reasoning proofs. We implement customized Prolog-based meta-interpreters to achieve two goals: (i) To produce reasoning proofs; (ii) To adopt better search algorithms other than DFS.

270 **3.2.1 Reasoning Tracer**
 271 We implement a Prolog meta-interpreter to show
 272 the reasoning proofs:

```

273 1 % Define the operator for proofs
274 2 :- op(750, xfy, =>).
275 3
276 4 % Proof tree generation
277 5 mi_tree(true, true).
278 6 mi_tree((A,B), (TA,TB)) :-
279 7     mi_tree(A, TA),
280 8     mi_tree(B, TB).
281 9 mi_tree(G, builtin(G)) :-
282 10    predicate_property(G, built-in
283    ),
284 11    !,
285 12    call(G).
286 13 mi_tree(g(G), TBody => G) :-
287 14    mi_clause(G, Body),
288 15    mi_tree(Body, TBody).
```

289 We showcase how a reasoning trace is induced
 290 using the example below. Given a knowledge base
 291 like:

```

292 1 parent_of(X, Y) :- mother_of(X, Y).
293 2 parent_of(X, Y) :- father_of(X, Y).
294 3 grandparent_of(X, Y) :-
295 4     parent_of(X, Z), parent_of(Z, Y).
296 5 mother_of(morty, beth).
297 6 father_of(beth, rick).
```

298 and a query:

```

299 1 ?- mi_tree(g(grandparent_of(morty, Who
300    )), Proof).
```

301 the output would be

```

302 1 Who=rick,
303 2 Proof=((((true=>mother_of(morty, beth))
304    =>parent_of(morty, beth), (true=>
305    father_of(beth, rick))=>parent_of(
306    beth, rick))=>grandparent_of(morty,
307    rick)).
```

308 The output proof is ensured to be causal and reli-
 309 able since a symbolic approach generates it.

3.2.2 Search Strategy

311 The default search strategy of Prolog is DFS, which
 312 may lead to infinite loops. For example, given the
 313 knowledge base:

```

314 1 parent_of(X, Y) :- offspring_of(Y, X).
315 2 offspring_of(X, Y) :- parent_of(Y, X).
316 3 parent_of(X, Y) :- mother_of(X, Y).
317 4 parent_of(X, Y) :- father_of(X, Y).
318 5 mother_of(jack, anna).
```

319 and a query ?- parent_of(jack, Who). , the
 320 backtracking process would repeat over the first
 321 two lines without resorting to other lines due to
 322 DFS. To address this issue, we adopt Iterative Deep-
 323 ening Search (IDS), in which the backtracking pro-
 324 cess performs a series of depth-limited searches,
 325 each with an increasing depth limit. This leverages

the strengths of both Breadth-First Search (BFS)
 and DFS. Our Prolog meta-interpreter for IDS is
 implemented as:

```

329 1 % Depth-limited meta-interpreter with
330   proof tree generation
331 2 mi_limit(true, true, N, N).
332 3 mi_limit((A,B), (TA,TB), N0, N) :-
333 4     mi_limit(A, TA, N0, N1),
334 5     mi_limit(B, TB, N1, N).
335 6 mi_limit(g(G), TBody => G, N0, N) :-
336 7     N0 #> 0,
337 8     N1 #= N0 - 1,
338 9     mi_clause(G, Body),
339 10    mi_limit(Body, TBody, N1, N).
340 11
341 12 % Iterative deepening with proof tree
342   generation
343 13 mi_id(Goal, Proof) :-
344 14     length(_, N),
345 15     mi_limit(Goal, Proof, N, _).
346 16
347 17 % Iterative deepening with maximum
348   depth with proof tree generation
349 18 mi_id_limit(Goal, Proof, MaxDepth) :-
350 19     between(1, MaxDepth, N),
351 20     mi_limit(Goal, Proof, N, _).
```

In practice, other search strategies can also be
 implemented according to the nature of the target
 problems, such as Uniform-Cost Search (UCS) and
 Beam Search.

4 Experiments

We briefly introduce our experimental settings in
 §4.1 and show the experiment results in §4.2.

4.1 Experimental Settings

We present our experimental settings in this section,
 including our implementation details of the two
 components (§4.1.1), a brief introduction of the
 adopted datasets (§4.1.2) and the baselines (§4.1.4).

4.1.1 Implementation

SymGen We adopt the Code-LLaMA (Rozière
 et al., 2023) family as the base LLMs to trans-
 late natural languages into Prolog representations.
 Our prompting paradigm is in a pure few-shot in-
 context-learning (ICL) prompting style, without
 detailed human-written instructions. Each ICL
 demonstration comprises a question and a piece
 of Prolog code.

SymInfer We adopt SWI-Prolog (Wielemaker
 et al., 2012) and PySwip¹ packages to implement
 the symbolic inference engine. We set the maxi-
 mum depth to be 20 for Iterative Deepening Search
 and the number of generated reasoning paths to 20.

¹<https://github.com/yuce/pyswip>

4.1.2 Datasets

We evaluate CARING on three popular complex reasoning datasets, including two logical reasoning datasets (ProofWriter (Tafjord et al., 2021) and PrOntoQA (Saparov and He, 2023)) and one arithmetic dataset (GSM8K (Cobbe et al., 2021; Ribeiro et al., 2023)).

ProofWriter ProofWriter (Tafjord et al., 2021) is a commonly-used logical reasoning dataset. It contains many small rulebases of facts and rules, expressed in English. Each rulebase has a set of questions (English statements) that can either be proven true or false using proofs of various depths, or the answer is “Unknown” (in open-world setting, OWA). The proofs can naturally be represented as directed acyclic graphs (DAGs). The dataset is divided into several sub-sets according to maximum proof depth, namely $\{0, \leq 1, \leq 2, \leq 3, \leq 5\}$. We follow previous work (Pan et al., 2023) to use a 600-instance subset sampled from the most difficult depth-5 test set. We also report additional results on the full depth-5 test set in Appendix §A.1.1.

PrOntoQA PrOntoQA (Saparov and He, 2023) is a synthetic question answering dataset designed for diagnosing the logical reasoning ability of LLMs. Each example aims to validate the feasibility of a statement given a context. We report results on two subsets so we can compare CARING with previous methods. As for the results reported in Table 4, we follow Pan et al. (2023) to adopt the most difficult depth-5 *fictional characters* sub-set, which contains 500 statement-context pairs. As for the results in Table 3, we use the subset adopted by Hao et al. (2023). Similar to ProofWriter, the proofs provided by the dataset can be naturally represented as DAGs.

GSM8K GSM8K (Cobbe et al., 2021) is a multi-step arithmetic reasoning dataset composed of high-quality grade school math word problems. The original GSM8K dataset contains reasoning explanations written in natural language, which raises difficulties in evaluating intermediate steps automatically. Recently, Ribeiro et al. (2023) released a subset that contains 270 questions annotated with structured reasoning proofs in the format of DAGs. We adopt this subset to enable the evaluation of reasoning proofs.

	Method	Acc (%)	Proof Sim (%)		
			All	Correct	
GPT-4*	CoT	67.41	–	–	
	ToT	70.33	–	–	
	CR	71.67	–	–	
	DetermLR	79.17	–	–	
	Logic-LM	79.66	–	–	
Code-LLaMA	7B	CoT	46.33	9.69	14.95
		Ours	91.00	72.91	84.39
	13B	CoT	46.50	15.69	25.86
		Ours	95.67	80.65	86.00
	34B	CoT	52.00	15.76	27.74
		Ours	96.50	81.02	86.12

Table 2: Results on the subset of ProofWriter adopted by Pan et al. (2023). The default setting is 2-shot. “All”: on all instances. “Correct”: on correctly-predicted instances. *All GPT-4 numbers are from Sun et al. (2023).

4.1.3 Evaluation Metrics

Ribeiro et al. (2023) proposed two novel metrics to evaluate the quality of the generated reasoning proofs in addition to the prevalent answer accuracy metric. Similar to them, we adopt the following metrics to evaluate both the answers and the generated reasoning proofs.

Answer Accuracy Answer accuracy measures a model’s ability to predict the correct answer. A prediction is deemed correct if it is (i) the same as the gold option for multi-choice problems and (ii) the same integer as the gold answer for arithmetic reasoning problems. This metric is the upper bound for other metrics since a reasoning graph would be marked as incorrect without evaluation if the answer is marked as incorrect. We report this metric for all datasets.

Reasoning Proof Similarity As shown in Figure 2, the problems that we are interested in naturally compose reasoning proofs in the format of directed acyclic graphs (DAGs). Reasoning proof similarity $\text{sim}(\mathcal{G}_g, \mathcal{G}_p)$ measures the graph similarity between the gold and the predicted reasoning graphs. We follow Ribeiro et al. (2023) to adopt the graph edit distance function $\delta(\mathcal{G}_g, \mathcal{G}_p)$. This function quantifies the graph edit distance by determining the minimum number of operations required over nodes and edges to transform one graph into the other, thereby enabling a comparison of \mathcal{G}_g and \mathcal{G}_p based on their structural similarities. The reasoning graph similarity is normalized to $[0, 1]$

Base LLM	#Param	#Shot	Method	Acc (%)	Proof Acc (%)	
					All	Correct
LLaMA-1*	33B	8-shot	CoT	87.8	64.8	–
			RAP	94.2	78.8	–
Code-LLaMA	13B	2-shot	CoT	80.2	52.4	53.4
			Ours	99.0	98.2	99.2

Table 3: Results on the PrOntoQA subset that was adopted by RAP (Hao et al., 2023) for comparison with their method. The results marked with * are from their paper.

	Method	Acc (%)	Proof Acc (%)		
		All	Correct	Correct	
GPT4*	CoT	98.8	–	–	
	Logic-LM	83.2	–	–	
Code-LLaMA	7B	CoT	52.0	24.8	28.5
		Ours	98.8	98.4	99.6
	13B	CoT	61.0	32.2	35.9
		Ours	99.4	98.8	99.4
	34B	CoT	82.8	41.0	41.0
		Ours	100.0	100.0	100.0

Table 4: Results on the depth-5 subset of PrOntoQA. The default setting is 2-shot. Results marked with * are GPT-4 results reported by Logic-LM (Pan et al., 2023).

as:

$$\text{sim}(\mathcal{G}_p, \mathcal{G}_g) = 1 - \frac{\delta(\mathcal{G}_p, \mathcal{G}_g)}{\max\{|N_p| + |E_p|, |N_g| + |E_g|\}} \quad (1)$$

where $|N_p|$ and $|E_p|$ denote the count of nodes and edges, respectively, within the predicted reasoning graph. A similar notation applies to $|N_g|$ and $|E_g|$, which represent the number of nodes and edges in the gold graph. Note that the reasoning graph similarity is set to zero if the predicted answer is incorrect. We report this metric for ProofWriter and GSM8K.

Reasoning Proof Accuracy This metric evaluates the exact match between the gold and the predicted reasoning proofs in terms of both reasoning graph structures and textual contents². The reasoning proof accuracy is either 1 or 0 for a single instance, making it a discrete version of reasoning proof similarity. Since this metric requires the dataset to have structured content to enable automatic evaluation, we can only apply it to PrOntoQA, which is specifically designed for easy parsing of the proofs.

²Note that our implementation here is simpler than that of Ribeiro et al. (2023) because we only apply this metric to PrOntoQA, which is easy to get evaluated.

	Method	Acc	Proof Sim (%)	
			All	Correct
7B	CoT	13.70	4.99	36.39
	Ours	12.22	6.57	53.72
13B	CoT	15.56	5.76	37.03
	Ours	21.48	11.66	54.26
34B	CoT	35.19	13.04	37.07
	Ours	42.22	22.91	54.25

Table 5: Results on GSM8K. The default setting is 5-shot.

4.1.4 Baselines

All baselines prompt the LLMs with few-shot in-context-learning (ICL) demonstrations.

Chain-of-Thought (CoT) CoT prompting (Wei et al., 2022) prompts LLMs with ICL demonstrations that contain both intermediate reasoning steps and answers. It serves as a popular and strong baseline for prompting LLMs to solve problems.

Logic-LM Logic-LM (Pan et al., 2023) is a neuro-symbolic method that adopts symbolic solvers for logical reasoning problems. The main difference between Logic-LM and CARING is: Logic-LM adopts various solvers for multiple datasets and only focuses on answer accuracy, while CARING universally uses one solver (i.e., SWI-Prolog); and more importantly, CARING showcases how SWI-Prolog interpreters can be customized to generate intermediate reasoning proofs and to adopt various search (i.e., problem-solving) strategies.

Search-based Methods We include several search-based methods as baselines. We directly adopt the released results in their papers, since it is too time-consuming to implement these methods on our own. For ProofWriter, we compare our method with GPT-4 based Tree-of-Thoughts (ToT; (Yao et al., 2023)), Cumulative Reasoning (CR; (Zhang et al., 2023)), and DetermLR (Sun

et al., 2023). We cannot make comparisons on reasoning proofs because these methods only reported reasoning accuracy. For PrOntoQA, we compare our method with RAP (Hao et al., 2023), in terms of both reasoning accuracy and reasoning proof accuracy.

4.2 Main Results

Tables 2, 3, 4 and 5 show the experimental results on our adopted datasets.

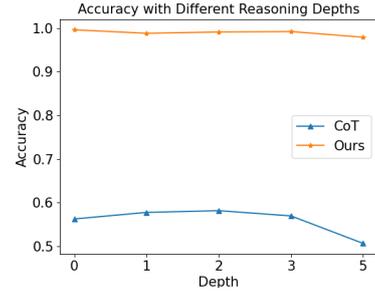
ProofWriter The results on ProofWriter are presented in Table 2. CARING demonstrates notable improvements over existing baselines, particularly in terms of reasoning proof similarity. Utilizing Code-LLaMA-34B, CARING achieves a remarkable answer accuracy of 96.50% and a reasoning proof similarity of 81.02%, significantly surpassing the most powerful method using GPT-4 that obtains an accuracy of 79.66%.

PrOntoQA The results on PrOntoQA are presented in Tables 3 and 4. CARING achieves almost full accuracy with the 13B model. Comparing with RAP, CARING obtains better results in terms of both answer accuracy and proof accuracy even using a smaller base LLM and fewer ICL demonstrations. CARING also outperforms CoT and LogicLM that use GPT-4.

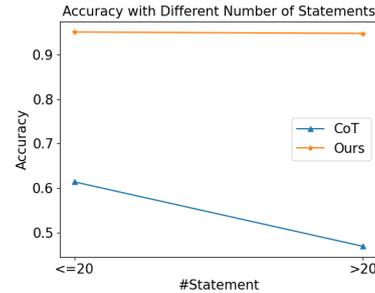
GSM8K The results on GSM8K are presented in Table 5. This dataset is more challenging than the previous two logical reasoning datasets for CARING, since it is generally believed that symbolic languages are restricted by their limited expressiveness and cannot properly handle the ambiguity in real-world human languages. Surprisingly, with the 34B model, CARING outperforms the strong CoT baseline by a large margin and almost doubles the reasoning proof similarity (22.91% vs. 13.04%). We attribute such improvements to increasingly powerful LLMs, which can correctly translate ambiguous human languages into formal symbolic representations.

4.3 When Reasoning Becomes More Complex

A key difficulty confronted by reasoning systems is the rapid expansion of possible states as the reasoning process becomes more complex, such as when additional statements are considered or the depth of inference is greater. To investigate how our method handles more complex reasoning problems, we conduct experiments under two controlled settings: (1)



(a) Answer accuracy with different reasoning depths.



(b) Answer accuracy with different number of statements.

Figure 4: Answer accuracy when reasoning problems become more complex.

#Depth \uparrow : How does the answer accuracy change with #Depth being ≤ 0 , ≤ 1 , ≤ 2 , ≤ 3 and ≤ 5 , respectively; (2) **#Statements** \uparrow : How does the answer accuracy change with #Statements being ≤ 20 and > 20 , respectively.

As shown in Figures 4a and 4b, with increasing levels of reasoning intricacy, the answer accuracy of CARING remains steady. In contrast, CoT sees significant decreases in answer accuracy under both settings. This verifies the robustness of CARING against complex reasoning.

5 Conclusion

This paper presents a framework to address the erroneous reasoning proof problem of LLM-based reasoning systems. Specifically, we develop a neuro-symbolic method called CARING, which produces high-quality reasoning proofs for complex reasoning problems. By implementing customized meta-interpreters for executing Prolog representations and putting LLMs under quarantine during the reasoning phase, CARING ensures the reasoning proofs to be causal and reliable. We conduct experiments on two logical reasoning datasets and one arithmetic reasoning dataset. Experimental results demonstrate our method achieves significant improvements with both final answers and intermediate reasoning proofs. Further analysis indicates CARING remains robust when the reasoning problems become more complex.

581 **Limitations**

582 We observe two limitations regarding our frame-
583 work:

- 584 • The generalization ability of CARING is re-
585 stricted by the expressiveness of the concern-
586 ing symbolic representations. In this paper,
587 we showcase a Prolog-based implementation.
588 Other symbolic representations could be ex-
589 plored to generalize CARING to more reason-
590 ing tasks.
- 591 • CARING requires powerful LLMs as sym-
592 bolic representation generators, which is sug-
593 gested by the results on GSM8K. This depen-
594 dence might prevent it from being applied to
595 productions.

596 **References**

597 Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and
598 Dan Klein. 2016. [Neural module networks](#). In *2016*
599 *IEEE Conference on Computer Vision and Pattern*
600 *Recognition, CVPR 2016, Las Vegas, NV, USA, June*
601 *27-30, 2016*, pages 39–48. IEEE Computer Society.

602 Maciej Besta, Nils Blach, Ales Kubicek, Robert Ger-
603 stenberger, Lukas Gianinazzi, Joanna Gajda, Tomasz
604 Lehmann, Michal Podstawski, Hubert Niewiadom-
605 ski, Piotr Nyczyk, and Torsten Hoeffler. 2023. [Graph](#)
606 [of thoughts: Solving elaborate problems with large](#)
607 [language models](#).

608 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,
609 Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias
610 Plappert, Jerry Tworek, Jacob Hilton, Reiichiro
611 Nakano, Christopher Hesse, and John Schulman.
612 2021. [Training verifiers to solve math word prob-](#)
613 [lems](#). *ArXiv preprint*, abs/2110.14168.

614 Alain Colmerauer and Philippe Roussel. 1996. *The*
615 *Birth of Prolog*, page 331–367. Association for Com-
616 puting Machinery, New York, NY, USA.

617 Antonia Creswell, Murray Shanahan, and Irina Higgins.
618 2023. [Selection-inference: Exploiting large language](#)
619 [models for interpretable logical reasoning](#). In *The*
620 *Eleventh International Conference on Learning Rep-*
621 *resentations*.

622 Nitish Gupta, Kevin Lin, Dan Roth, Sameer Singh, and
623 Matt Gardner. 2020. [Neural module networks for rea-](#)
624 [soning over text](#). In *8th International Conference on*
625 *Learning Representations, ICLR 2020, Addis Ababa,*
626 *Ethiopia, April 26-30, 2020*. OpenReview.net.

627 Shibo Hao, Yi Gu, Haodi Ma, Joshua Hong, Zhen
628 Wang, Daisy Wang, and Zhiting Hu. 2023. [Reason-](#)
629 [ing with language model is planning with world](#)
630 [model](#). In *Proceedings of the 2023 Conference on*

Empirical Methods in Natural Language Processing,
pages 8154–8173, Singapore. Association for Com-
putational Linguistics. 631
632
633

Drew A. Hudson and Christopher D. Manning. 2019. [Learning by abstraction: The neural state machine](#).
In *Advances in Neural Information Processing Sys-*
tems 32: Annual Conference on Neural Information
Processing Systems 2019, NeurIPS 2019, December
8-14, 2019, Vancouver, BC, Canada, pages 5901–
5914. 634
635
636
637
638
639
640

Mehran Kazemi, Najoung Kim, Deepti Bhatia, Xin
Xu, and Deepak Ramachandran. 2023. [LAMBADA:](#)
[Backward chaining for automated reasoning in nat-](#)
[ural language](#). In *Proceedings of the 61st Annual*
Meeting of the Association for Computational Lin-
guistics (Volume 1: Long Papers), pages 6547–6568,
Toronto, Canada. Association for Computational Lin-
guistics. 641
642
643
644
645
646
647
648

Robert Kowalski. 1979. [Algorithm = logic + control](#).
Commun. ACM, 22(7):424–436. 649
650

Jieyi Long. 2023. [Large language model guided tree-of-](#)
[thought](#). 651
652

Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang,
Delip Rao, Eric Wong, Marianna Apidianaki, and
Chris Callison-Burch. 2023. [Faithful chain-of-](#)
[thought reasoning](#). *ArXiv preprint*, abs/2301.13379. 653
654
655
656

Arvind Neelakantan, Quoc V. Le, Martín Abadi, An-
drew McCallum, and Dario Amodei. 2017. [Learning](#)
[a natural language interface with neural program-](#)
[mer](#). In *5th International Conference on Learning*
Representations, ICLR 2017, Toulon, France, April
24-26, 2017, Conference Track Proceedings. Open-
Review.net. 657
658
659
660
661
662
663

Maxwell Nye, Michael Henry Tessler, Joshua B. Tenen-
baum, and Brenden M. Lake. 2021. [Improving co-](#)
[herence and consistency in neural sequence models](#)
[with dual-system, neuro-symbolic reasoning](#). 664
665
666
667

OpenAI. 2023. [Gpt-4 technical report](#). 668

Liangming Pan, Alon Albalak, Xinyi Wang, and
William Yang Wang. 2023. [Logic-LM: empower-](#)
[ing large language models with symbolic solvers for](#)
[faithful logical reasoning](#). In *Findings of the 2023*
Conference on Empirical Methods in Natural Lan-
guage Processing (Findings of EMNLP), Singapore. 669
670
671
672
673
674

Danilo Neves Ribeiro, Shen Wang, Xiaofei Ma,
Henghui Zhu, Rui Dong, Deguang Kong, Juli-
ette Burger, Anjelica Ramos, zhiheng huang,
William Yang Wang, George Karypis, Bing Xiang,
and Dan Roth. 2023. [STREET: A MULTI-TASK](#)
[STRUCTURED REASONING AND EXPLANA-](#)
[TION BENCHMARK](#). In *The Eleventh Interna-*
tional Conference on Learning Representations. 675
676
677
678
679
680
681
682

Baptiste Rozière, Jonas Gehring, Fabian Gloeckle,
Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi
Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom 683
684
685

		Acc	Proof Sim	
			All	Correct
7B	Direct	41.78	–	–
	Direct (3-Shot)	43.32	–	–
	CoT	40.95	11.27	17.20
	CoT (3-shot)	42.58	11.52	21.58
	Ours	92.43	75.85	86.68
	<hr/>			
13B	Direct	43.44	–	–
	Direct (3-shot)	44.31	–	–
	CoT	45.88	16.16	27.32
	CoT (3-shot)	54.70	23.18	32.48
	Ours	96.16	80.74	86.34
<hr/>				
34B	Direct	44.00	–	–
	Direct (3-shot)	45.93	–	–
	CoT	52.32	15.08	26.30
	CoT (3-shot)	56.50	24.12	34.61
	Ours	98.11	83.17	85.65

Table 6: Results on ProofWriter. “All” and “Correct” refer to “on all instances” and “on correctly-predicted instances”, respectively. “Proof Sim” refers to “Proof Graph Similarity” while “Proof EM” means “Proof Graph Exact Match”. The default setting is 2-shot. We additionally conduct 3-shot experiments for baselines to include all types of labels in the in-context demonstrations because this dataset contains three labels: {true, false, uncertain}. We do not conduct 3-shot experiments for our method because it is not sensitive to the number of labels due to its reasoning-by-execution nature.

A Appendix

A.1 Additional Results

A.1.1 Results on ProofWriter

Table 6 shows the results from our implementation on the depth-5 test set of ProofWriter.