

# 000 001 002 003 004 005 $\pi$ -CoT: PROLOG-INITIALIZED CHAIN-OF-THOUGHT 006 PROMPTING FOR MULTI-HOP QUESTION-ANSWERING 007 008 009

010 **Anonymous authors**  
011

012 Paper under double-blind review  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022

## 023 ABSTRACT 024

025 Chain-of-Thought (CoT) prompting significantly enhances large language models'  
026 (LLMs) problem-solving capabilities, but still struggles with complex multi-hop  
027 questions, often falling into circular reasoning patterns or deviating from the logical  
028 path entirely. This limitation is particularly acute in retrieval-augmented generation  
029 (RAG) settings, where obtaining the right context is critical. We introduce  
030 **Prolog-Initialized Chain-of-Thought** ( $\pi$ -CoT), a novel prompting strategy that  
031 combines logic programming's structural rigor with language models' flexibility.  
032  $\pi$ -CoT reformulates multi-hop questions into Prolog queries decomposed as  
033 single-hop sub-queries. These are resolved sequentially, producing intermediate  
034 artifacts, with which we initialize the subsequent CoT reasoning procedure.  
035 Extensive experiments demonstrate that  $\pi$ -CoT significantly outperforms standard  
036 RAG and in-context CoT on multi-hop question-answering benchmarks.  
037

## 038 1 INTRODUCTION 039

040 Chain-of-thought (CoT) reasoning has emerged as a powerful paradigm for enhancing the  
041 problem-solving capabilities of large language models, substantially improving performance on  
042 arithmetic, commonsense, and symbolic reasoning tasks (Wei et al., 2022; Kojima et al., 2022). By  
043 encouraging models to articulate their reasoning process through intermediate steps, CoT enables  
044 more systematic and interpretable problem-solving approaches (Zhang et al., 2022; Wang et al., 2022).  
045 However, as the complexity of reasoning tasks increases—particularly in multi-hop scenarios where  
046 multiple interconnected inferences must be made—CoT systems have been observed to generalize  
047 poorly (Dziri et al., 2023) and become trapped in circular reasoning patterns (Lo et al., 2023; Yao  
048 et al., 2023).

049 This limitation becomes especially pronounced in retrieval-augmented generation (RAG) systems,  
050 where CoT excels at single-hop questions that require straightforward document retrieval and  
051 reasoning, but struggles significantly with multi-hop queries that demand the integration of  
052 information across multiple sources and reasoning steps (Asai et al., 2023). The fundamental  
053 challenge lies in CoT's inherent trade-off: while its flexibility allows for creative and adaptive  
054 reasoning, this same flexibility can lead to unstructured exploration that fails to maintain logical  
055 consistency across complex reasoning chains.

056 Recent work has explored decomposition-based approaches that break multi-hop questions into  
057 manageable single-hop questions (Khot et al., 2023; Zhou et al., 2023; Min et al., 2019).  
058 However, even with decomposition, critical gaps remain: models struggle to generate high-quality  
059 decompositions without supervision (Patel et al., 2022; Wolfson et al., 2020), fail at reliable fact  
060 composition across steps (Press et al., 2023), and lose track of intermediate state in long reasoning  
061 chains (Yen et al., 2024; Liu et al., 2024). These failures suggest the need for a more principled  
062 reasoning framework that can enforce structure while maintaining flexibility.

063 In contrast to natural-language reasoning pioneered by CoT, the structured reasoning paradigm has  
064 been extensively studied for decades in artificial intelligence and logic programming (Kowalski and  
065 Smoliar, 1982; Russell et al., 1995; McCarthy et al., 1960; Simon and Newell, 1971). Prolog, a  
066 declarative programming language explicitly designed for structured reasoning tasks, exemplifies  
067 this approach through its systematic query resolution mechanisms and logical rule-based inference  
068 (Robinson, 1965; Kowalski and Smoliar, 1982; Clocksin and Mellish, 2003). While Prolog's rigid  
069

054 structure ensures logical consistency, it lacks the flexibility to handle ambiguous natural language,  
 055 cannot easily incorporate unstructured text from documents, and requires precise logical formulations  
 056 that may not capture the nuanced reasoning needed for real-world questions.  
 057

058 Recognizing that Prolog and CoT possess complementary strengths and weaknesses, we introduce  
 059 **Prolog-Initialized Chain-of-Thought** ( $\pi$ -CoT), a novel prompting strategy that combines the  
 060 structural rigor of logic programming with the contextual flexibility of natural language reasoning.  
 061 Our approach begins by algorithmically reformulating complex multi-hop reasoning questions into  
 062 equivalent Prolog queries, where each query is deliberately decomposed into a sequence of single-hop  
 063 sub-queries. These sub-queries are then resolved systematically: each is translated into natural  
 064 language and posed to a RAG (Lewis et al., 2020; Gao et al., 2023) or in-context CoT system, which  
 065 retrieves relevant documents and generates answers. The resulting answers are translated back into  
 066 Prolog facts and incorporated into the evolving knowledge base.  
 067

068 The key insight underlying  $\pi$ -CoT is that by structuring the reasoning process through Prolog’s  
 069 query resolution mechanism, we ensure that the retrieved context remains highly relevant. Rather  
 070 than allowing the model to freely explore the reasoning space, potentially losing track of relevant  
 071 information or pursuing irrelevant tangents (Yao et al., 2023; Dziri et al., 2023), our approach  
 072 maintains a structured trajectory that systematically builds toward the final answer.  
 073

074 At the completion of the Prolog resolution process, we concatenate the original question, all retrieved  
 075 documents, and the structured Prolog derivation to create a comprehensive context that initializes  
 076 the final CoT reasoning step. Because most of the heavy-lifting already happens in the initialization  
 077 process, the final CoT reasoning is far simpler and more successful.  
 078

079 Through extensive experimental evaluation, we demonstrate that  $\pi$ -CoT is on par or better than  
 080 traditional RAG and in-context systems on multi-hop question-answering (QA) benchmarks, including  
 081 HotpotQA, 2WikiMultiHopQA, MuSiQue, and PhantomWiki. Our results suggest that the principled  
 082 integration of symbolic reasoning structures with neural language models offers a promising direction  
 083 for developing more reliable and interpretable reasoning systems.  
 084

## 085 2 RELATED WORKS

086 **Decomposition for multi-hop question-answering.** Breaking down a complex problem into smaller,  
 087 manageable parts is a common technique in LLM prompting (Zhou et al., 2023; Khot et al., 2023;  
 088 Wei et al., 2022). For open-domain QA, Press et al. (2023) prompt the model to generate follow-up  
 089 questions, and Trivedi et al. (2023) take each new sentence in a CoT as input to the retriever.  
 090 Importantly, the language model decomposes the question in natural-language steps. Recent works  
 091 have also explored the use of *explicit plans*, usually in the form of Python programs (Surís et al.,  
 092 2023; Khattab et al., 2022). While we do not provide a direct comparison due to different model sizes  
 093 and/or retrieval setups, Monte Carlo Tree Search (Tran et al., 2024), test-time scaling (Wang et al.,  
 094 2025), and reinforcement learning methods (Li et al., 2025; Jin et al., 2025a; Song et al., 2025) are  
 095 emerging as promising approaches to open-domain QA.  
 096

097 **Improving language model reasoning with Prolog.** Many prior works generate Prolog from natural  
 098 language to improve arithmetic reasoning or multi-hop question-answering. Wu and Liu (2025);  
 099 Vakharia et al. (2024); Borazjanizadeh and Piantadosi (2024) translate questions into Prolog, then  
 100 query a knowledge base. However, they assume the knowledge base has already been populated  
 101 with facts. Therefore, they do not address retrieval from unstructured documents. Tan et al. (2024);  
 102 Yang et al. (2024) utilize Prolog as a source of supervised training signals for math and logical  
 103 reasoning. Weber et al. (2019) propose a weak unification strategy in Prolog based on semantic  
 104 similarity. However, their approach requires training and uses pre-defined predicates extracted from  
 105 the training text. The method closest to our work is that of Chen et al. (2019). They train an LSTM  
 106 “programmer” to generate programs, which are executed using a BERT-based “reader” to produce  
 107 answers. In this work, we contribute a training-free strategy and demonstrate its effectiveness with  
 108 recent LLMs like Llama-3.3-70B-Instruct and Deepseek-R1-Distill-Qwen-32B, using both sparse  
 109 and dense retrievers. The results of Chen et al. (2019) are also limited by the strictness of a pure  
 110 Prolog execution.  
 111

112 **Fact extraction and summarization.** Extracting knowledge graph triples from unstructured text  
 113 is a classical problem in NLP, also known as open information extraction (OpenIE) (Angeli et al.,  
 114

108 2015; Pei et al., 2023; Zhou et al., 2022). To enhance conventional retrieval techniques, LightRAG  
 109 (Guo et al., 2024) and HippoRAG (Jimenez Gutierrez et al., 2024; Gutiérrez et al., 2025) demonstrate  
 110 the effectiveness of fact extraction and GraphRAG (Edge et al., 2024) and RAPTOR (Sarthi et al.,  
 111 2024) propose methods for clustering and summarization. Among these methods, HippoRAG 2  
 112 from Gutiérrez et al. (2025) performs the best on HotpotQA, 2WikiMultiHopQA, and MuSiQue. We  
 113 provide a direct comparison of our method to HippoRAG 2 in the results section below.  
 114

### 115 3 PRELIMINARIES

117 Dating back to the 1970s, Prolog<sup>1</sup> is a powerful way to represent factual knowledge and perform  
 118 logical inference (Colmerauer and Roussel, 1996; Russell et al., 1995; Sterling and Shapiro, 1994).  
 119 Solutions in Prolog are verifiable and compositional, making it particularly well-suited for multi-hop  
 120 question answering where intermediate steps must be chained reliably.  
 121

122 **Representing factual knowledge.** Consider the the following English sentence about *Harry Potter*  
 123 (Rowling, 1997):

124 “Lockhart is a Defense against the Dark Arts teacher at Hogwarts.”

125 This statement can be represented as the following Prolog **fact**:

126 `DADA_teacher ("Lockhart", "Hogwarts").`

127 In Prolog, `DADA_teacher` is a **predicate** and “Lockhart” and “Hogwarts” are **assignments**. A  
 128 large amount of real-world knowledge can be encoded in this structured, relational form. In fact, as  
 129 of early 2025, Wikidata contained over 1.65 billion such knowledge triples.<sup>2</sup> Throughout the paper,  
 130 we refer to a collection of Prolog facts as a **knowledge base**.  
 131

132 **Multi-hop questions as Prolog queries.** Prolog doesn’t just store facts—it also allows us to query  
 133 them. A **Prolog query** comprises one or more predicates with unassigned variables and asks whether  
 134 any satisfying variable assignment exists. Consider the following example:

135 “Who is the wife of the Defense against the Dark Arts teacher at Hogwarts  
 136 (1)  
 137 who is also a werewolf?”

138 One way of answering this question involves first identifying all the Defense against the Dark Arts  
 139 teachers at Hogwarts, then filtering for those who are also werewolves, and finally retrieving the wives  
 140 of the selected people. Each step depends on resolving the previous step(s), and the intermediate  
 141 space of possible answers can be large. (Note that J.K. Rowling’s books mention seven Defense  
 142 against the Dark Arts teachers at Hogwarts, so a language model would have to consider seven  
 143 separate reasoning paths to answer the question.) A concise way of expressing question (1) is the  
 144 following Prolog query:  
 145

146 `DADA_teacher(X, "Hogwarts"),`  
 147 `werewolf(X),`  
 148 `wife(X, Y),` (2)

149 with `Y` representing the answer. Translating questions in natural language to structured queries is a  
 150 classical problem in the community (Zelle and Mooney, 1996; Zettlemoyer and Collins, 2012).

151 A **solution** to a Prolog query is a set of variable assignments. Assuming we have a  
 152 pre-populated knowledge base, the solutions are obtained by executing the query. For example,  
 153 a knowledge base containing the facts `DADA_teacher("Lockhart", "Hogwarts")`,  
 154 `werewolf("Lupin")`, and `wife("Lupin", "Tonks")` yields the following solution to  
 155 query (2):  
 156

157 `X = "Lupin",`  
 158 `Y = "Tonks".`

159 <sup>1</sup>We specifically use the SWI-Prolog implementation (Wielemaker et al., 2012) of Prolog due to its rich  
 160 support of aggregation and if-then-else logic, which allows us to resolve questions like “How many Defense  
 161 against the Dark Arts teachers have been at Hogwarts?”.

<sup>2</sup><https://en.wikipedia.org/wiki/Wikidata>

162 4 METHOD:  $\pi$ -CoT — PROLOG-INITIALIZED CHAIN-OF-THOUGHT  
163

164 Prolog and LLMs are both powerful tools  
165 for reasoning, but they have complementary  
166 strengths and weaknesses. While Prolog  
167 provides a precise and verifiable framework  
168 for multi-hop reasoning, it assumes access to  
169 a *structured* database, which is often difficult  
170 to obtain in practice. In contrast, LLMs can  
171 adeptly retrieve and extract information from  
172 *unstructured* text, but cannot guarantee logical  
173 consistency across reasoning steps.

174 Motivated by these observations, we combine  
175 Prolog and LLMs in two ways: a) we initialize  
176 the CoT prompt with the execution trace of a  
177 Prolog query. This mitigates LLMs’ tendency  
178 to diverge from successful reasoning paths for  
179 CoT prompting.

180 b) we retrieve relevant information from  
181 available documents with the LLM and generate  
182 Prolog facts. This provides Prolog an indirect  
183 mechanism to interface natural-language  
184 documents and create a structured database.

185 **Prolog-Initialized CoT.** Our resulting  
186 workflow ( $\pi$ -CoT) is illustrated in Fig. 1.  
187 Given the question in natural language, we  
188 first prompt<sup>3</sup> the LLM to generate a structured  
189 **Prolog query:**

$$Q = (q_1, q_2, \dots, q_T).$$

190 Here,  $q_i$  is a **sub-query**. The query  $Q$  can be  
191 resolved step-by-step, one sub-query at a time.  
192 By design, the **answer** to the original question  
193 lies in one of these sub-queries, usually the last  
194 one. For example, the question in eq.(1) yields  
195 the Prolog query (2) with  $T = 3$  and answer  $Y$ .  
196

197  $\pi$ -CoT resolves the query step-by-step and logs  
198 the execution trace along the way. At step  $t$ ,  
199  $\pi$ -CoT stores the set of all possible solutions,  $S_t$ ,  
200 to the partially formed Prolog query  $(q_1, \dots, q_t)$ .  
201 Each element in  $S_t$  is a dictionary of key-value  
202 pairs, where the keys are the names of all  
203 variables (e.g.  $X, Y$ ) in the Prolog query and the  
204 values are valid assignments (e.g. Lupin).

205 4.1 SINGLE-STEP EXECUTION WITH SLICE  
206

207 Each step  $t$  follows a fixed procedure, which we refer to as **SLICE**<sup>4</sup>. Given  $S_{t-1}$  from the previous  
208 iteration, we resolve the next sub-query  $q_t$  and create  $S_t$ . The sub-query  $q_t$  can be of exactly two  
209 types: **Verification** (e.g. `werewolf(X)`), or **Extraction** (e.g. `wife(X, Y)`). Verification queries  
210 reduce the set of solutions in  $S_{t-1}$  to those compatible with the query (e.g. removing all teacher  
211 names assigned to  $X$  that are not werewolves). This scenario is shown in Fig. 2 for sub-query  $q_2$ .  
212 Extraction queries add new variables to  $S_t$  that satisfy the query constraint (e.g. adding variable  
213  $Y$  and assigning it the names of all teachers’ wives already assigned to  $X$ ). We resolve both query  
214 types with the help of the LLM and the available reference documents, while logging the retrieved

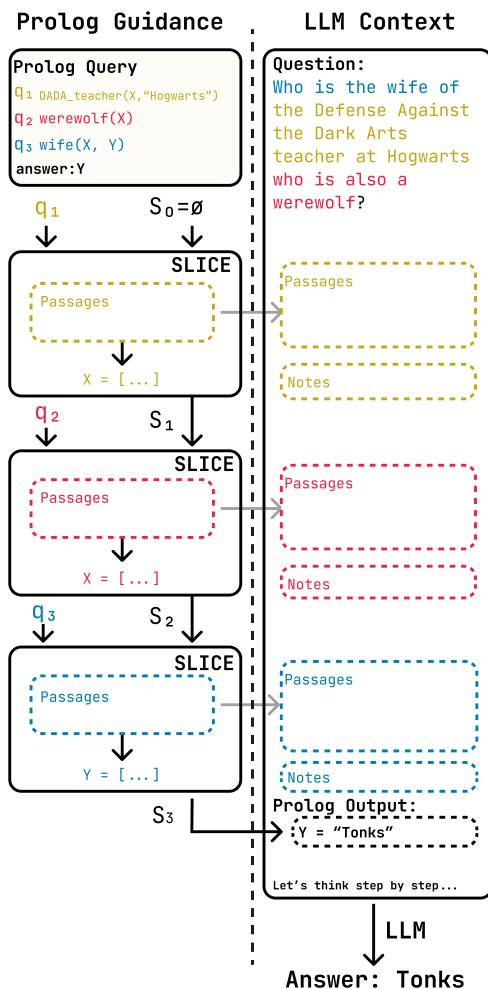


Figure 1: **Overview of  $\pi$ -CoT.** Left:  $\pi$ -CoT executes an LLM-generated Prolog query, using the **SLICE** module to resolve each sub-query  $q_t$ . Right:  $\pi$ -CoT uses the passages, notes, and (potentially) answer from the **SLICE** modules to initialize the CoT prompt for the final LLM call.

<sup>3</sup>We provide the query generation prompt in App. B.1.

<sup>4</sup>SLICE stands for Single-step Logical Inference with Contextual Evidence

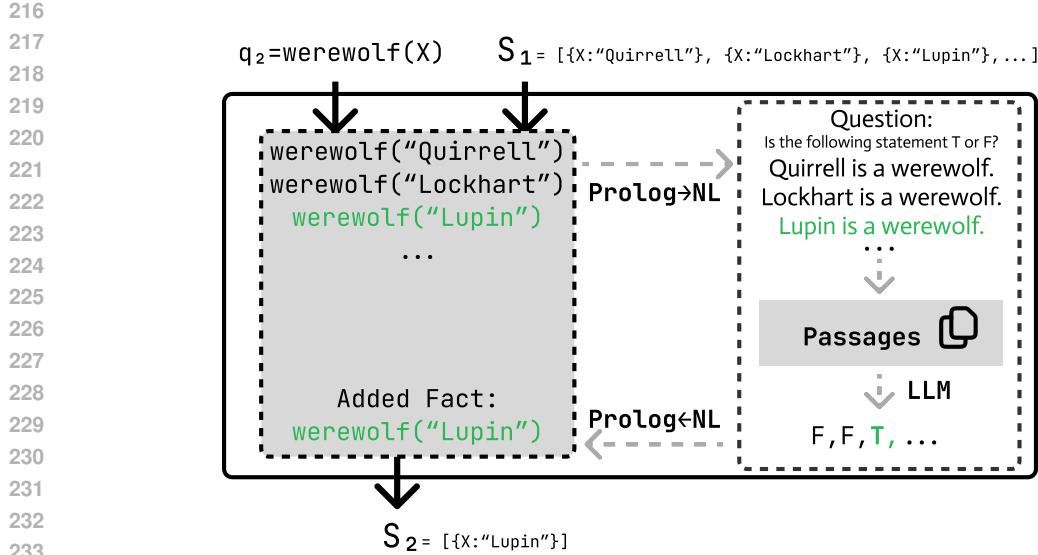


Figure 2: **SLICE module for fact verification in the RAG setting.** At step  $t = 2$ , the module takes in the previous state  $S_1$  containing variable assignments, the current sub-query  $q_2$ , and the corpus  $\mathcal{C}$  (not shown) as inputs and outputs  $S_2$ . Only the Prolog fact (in green) corresponding to a valid statement is added to a growing knowledge base.

document *passages* and Prolog facts (rephrased in natural language as *Notes*). The Prolog facts are stored in a Prolog knowledge base, allowing us to execute the query ( $q_1, \dots, q_t$ ) with an off-the-shelf Prolog interpreter and log its answer. We initialize  $S_0 = \emptyset$ , as the empty set. We provide additional details of the SLICE module in App. A.

#### 4.2 SLICE CHAINING

The inputs to the SLICE module are the sub-query  $q_t$ , the previous solutions  $S_{t-1}$ , and the document corpus  $\mathcal{C}$ . The output of the SLICE module is  $S_t$ . To derive the final solution, we chain the SLICE modules as follows:

$$S_t = \text{SLICE}(q_t, S_{t-1}, \mathcal{C}) \quad \text{for } t = 1, 2, \dots, T \quad \text{with } S_0 = \emptyset.$$

By design, the set of solutions is initially empty (i.e.,  $S_0 = \emptyset$ ). As the Prolog execution progresses—and more Prolog facts are collected— $S_t$  gets closer to the final solution. After all sub-queries are resolved,  $S_T$  is a set of solutions containing the final answer. For example, Fig. 1 shows  $S_3 = \{Y : "Tonks"\}$  as the final solution.

#### 4.3 COMBINING SYMBOLIC AND NATURAL LANGUAGE REASONING

In summary, the iterative process of SLICE chaining generates the following artifacts:

- a collection of retrieved **passages** from each SLICE execution<sup>5</sup>;
- a collection of Prolog facts, which we convert into natural language (i.e., “**notes**”); and
- the **answer** from Prolog execution.

The passages contain the necessary factual context to answer the original question. Since these passages may also contain superfluous information, the notes help to focus the model only on relevant information. Due to their iterative construction, these notes also serve to guide the LLM like breadcrumbs toward the final answer. When the final solution is non-empty, the LLM merely needs

<sup>5</sup>In the in-context setting, the retrieved passages and original passages are the same.

270 Table 1: **Accuracy of prompting-based methods on open-domain QA.** We report mean  $\pm$  1 standard  
 271 error for exact match (EM) and F1 score across 500 randomly chosen questions from each dataset.  
 272 Given a dataset and metric, we perform a repeated measures ANOVA followed by Tukey’s HSD with  
 273  $\alpha = 0.05$  to test significance of paired differences between methods. Bold indicates that no other  
 274 method performs significantly better.

296 297 298 299 300 301 302 303 304 305 306 307 308 309 310 311 312 313 314 315 316 317 318 319 320 321 322 323	296 297 298 299 300 301 302 303 304 305 306 307 308 309 310 311 312 313 314 315 316 317 318 319 320 321 322 323	302 303 304 305 306 307 308 309 310 311 312 313 314 315 316 317 318 319 320 321 322 323	302 303 304 305 306 307 308 309 310 311 312 313 314 315 316 317 318 319 320 321 322 323			
	Method	HotpotQA	2WikiMultiHopQA	MuSiQue		
	EM $\uparrow$	F1 $\uparrow$	EM $\uparrow$	F1 $\uparrow$	EM $\uparrow$	F1 $\uparrow$
Standard RAG	<b>38.8 <math>\pm</math> 2.2</b>	<b>52.6 <math>\pm</math> 2.0</b>	37.2 $\pm$ 2.2	40.4 $\pm$ 2.1	11.0 $\pm$ 1.4	18.1 $\pm$ 1.5
Self-Ask	19.2 $\pm$ 1.8	28.0 $\pm$ 1.8	15.8 $\pm$ 1.6	21.8 $\pm$ 1.7	5.6 $\pm$ 1.0	9.1 $\pm$ 1.2
IRCoT	<b>40.4 <math>\pm</math> 2.2</b>	<b>52.9 <math>\pm</math> 2.0</b>	32.4 $\pm$ 2.1	42.5 $\pm$ 2.0	<b>17.6 <math>\pm</math> 1.7</b>	<b>24.5 <math>\pm</math> 1.8</b>
$\pi$ -CoT (Ours)	<b>42.0 <math>\pm</math> 2.2</b>	<b>59.1 <math>\pm</math> 1.9</b>	<b>49.4 <math>\pm</math> 2.2</b>	<b>57.5 <math>\pm</math> 2.1</b>	<b>15.2 <math>\pm</math> 1.6</b>	<b>25.7 <math>\pm</math> 1.7</b>

284 Table 2: **Efficiency of prompting-based methods on open-domain QA.** We report mean  $\pm$  1  
 285 standard error number of BM25 queries, LLM calls, and total tokens for the same questions used in  
 286 Tab. 1. We further break down the total tokens into prompt and completion tokens in Tab. 7.

288 289 290 291 292 293 294 295	296 297 298 299 300 301 302 303 304 305 306 307 308 309 310 311 312 313 314 315 316 317 318 319 320 321 322 323	302 303 304 305 306 307 308 309 310 311 312 313 314 315 316 317 318 319 320 321 322 323	302 303 304 305 306 307 308 309 310 311 312 313 314 315 316 317 318 319 320 321 322 323						
	Method	HotpotQA	2WikiMultiHopQA	MuSiQue	BM25	LLM	Tokens	BM25	LLM
Standard RAG	1	1	$3.6 \times 10^3$	1	1	$3.7 \times 10^3$	1	1	$2.4 \times 10^3$
Self-Ask	3.36	3.36	$1.5 \times 10^4$	3.44	3.44	$1.6 \times 10^4$	3.29	3.29	$1.2 \times 10^4$
IRCoT	3.07	3.07	$6.2 \times 10^4$	3.47	3.47	$5.9 \times 10^4$	3.51	3.51	$4.5 \times 10^4$
$\pi$ -CoT (Ours)	2.82	4.82	$1.8 \times 10^4$	2.14	4.14	$2.2 \times 10^4$	3.42	5.42	$1.5 \times 10^4$

297 to return it. For the best results, we provide all three artifacts to the LLM and invoke chain-of-thought  
 298 reasoning to produce the final answer<sup>6</sup>.

## 5 MAIN RESULTS

303 We consider two settings: (1) **open-domain question-answering**, when the corpus is too large to fit  
 304 within the model’s context window, and (2) **in-context question-answering**, when the corpus does  
 305 fit within the model’s context window. The open-domain and in-context QA settings are also known  
 306 as the *fullwiki* and *distractor* settings in the literature (Yang et al., 2018). To overcome the context  
 307 limitations in the open-domain QA setting, we use retrieval augmented generation (RAG) to first  
 308 fetch relevant passages before generation. Since  $\pi$ -CoT is a prompting method, we require a strong  
 309 instruction-tuned model and employ the Llama-3.3-70B-Instruct model from Grattafiori et al. (2024).  
 310 We also provide results utilizing the Deepseek-R1-Distill-Qwen-32B model from (Guo et al., 2025)  
 311 in the in-context question-answering setting.

### 5.1 OPEN-DOMAIN QUESTION-ANSWERING

312 We evaluate  $\pi$ -CoT on three multi-hop QA datasets: (1) HotpotQA from Yang et al. (2018), (2)  
 313 2WikiMultiHopQA from Ho et al. (2020), and (3) MuSiQue from Trivedi et al. (2022). Since these  
 314 datasets are curated from Wikipedia, we can assess Prolog’s effectiveness in handling real-world  
 315 knowledge. For our retrieval setup, we use the preprocessed December 18, 2020 corpus from  
 316 FlashRAG (Jin et al., 2025b), which contains 20M chunks each of size 100 words, and use BM25  
 317 (Robertson et al., 2009) as our retriever. We provide supplementary experiment details in Sec. C.1.

318 Tab. 1 compares  $\pi$ -CoT to standard RAG and two multi-hop RAG baselines that also rely on  
 319 decomposition (via prompting) to handle multi-hop reasoning: (1) **Self-Ask** from Asai et al.

320  
321  
322  
323  
324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377  
378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
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  
429  
430  
431  
432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1139  
1130  
1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164

324  
 325 Table 3: **Comparison to the state-of-the-art OpenIE method.** We report exact match (EM) and F1  
 326 on the splits from [Gutiérrez et al. \(2025, Tab. 2\)](#). We use Llama-3.3-70B-Instruct for generation and  
 327 NV-Embed-v2 for embedding passages with  $k = 5$  per query.

328 329 330 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350 351 352 353 354 355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 377 Method	328 329 330 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350 351 352 353 354 355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 377 HotpotQA		328 329 330 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350 351 352 353 354 355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 377 2WikiMultiHopQA		328 329 330 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350 351 352 353 354 355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 377 MuSiQue	
	328 329 330 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350 351 352 353 354 355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 377 EM $\uparrow$	328 329 330 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350 351 352 353 354 355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 377 F1 $\uparrow$	328 329 330 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350 351 352 353 354 355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 377 EM $\uparrow$	328 329 330 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350 351 352 353 354 355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 377 F1 $\uparrow$		
328 329 330 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350 351 352 353 354 355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 377 Standard RAG	61.2	74.5	58.3	63.2	34.9	44.8
328 329 330 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350 351 352 353 354 355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 377 HippoRAG 2 (Gutiérrez et al., 2025)	<b>62.6</b>	75.3	65.5	72.0	37.6	49.5
328 329 330 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350 351 352 353 354 355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 377 π-CoT (Ours)	60.3	<b>76.8</b>	<b>71.1</b>	<b>79.6</b>	<b>38.5</b>	<b>56.2</b>

(2023) and (2) **IRCoT** from (Trivedi et al., 2023). Among all training-free<sup>7</sup> methods in the FlashRAG repository, IRCoT was the top-performing training-free method on HotpotQA and the second top-performing method on 2WikiMultiHopQA at the time of writing. On HotpotQA, standard RAG, IRCoT, and π-CoT are comparable in terms of accuracy, with neither method significantly outperforming the other two as determined by a Tukey’s HSD test with  $\alpha = 0.05$ . On 2WikiMultiHopQA, π-CoT significantly outperforms all other methods in terms of exact match and F1 score. On MuSiQue, IRCoT and π-CoT achieve comparable accuracy. Despite requiring only a single retriever call, standard RAG achieves surprisingly competitive accuracy on HotpotQA. Min et al. (2019); Chen and Durrett (2019) reveal that multi-hop reasoning is not required for many examples in HotpotQA, possibly explaining our findings. Notably, we find that π-CoT never performs worse than standard RAG, and does significantly better in the case of 2WikiMultiHopQA and MuSiQue.

To assess the efficiency of the methods in Tab. 1, we measure the number of BM25 queries, the number of LLM calls, and the total token usage. As shown in Tab. 2, π-CoT uses a similar number of BM25 queries as IRCoT and Self-Ask. We also observe π-CoT using more LLM calls on average than IRCoT and Self-Ask, which makes sense given that π-CoT must generate a Prolog query and perform the final chain-of-thought reasoning step. For the price of two extra LLM call, π-CoT enjoys lower total token usage. In particular, the use of Prolog allows intermediate steps to use separate, short contexts during execution, rather than a single, long context that grows with the number of steps. Our results in Tab. 2 reflect this fundamental difference.

Next, we compare π-CoT to **HippoRAG 2**, an OpenIE-augmented retrieval method that outperforms GraphRAG (Edge et al., 2024), RAPTOR (Sarthi et al., 2024), and LightRAG (Guo et al., 2024) on multi-hop QA. We follow the experimental setup of Gutiérrez et al. (2025, Tab. 2), which uses the NV-Embed-v2 embedding model of Lee et al. (2024) for retrieval. Instead of using full Wikipedia as the corpus, this experiment uses 9811 passages for HotpotQA, 6119 passages for 2WikiMultiHopQA, and 11656 passages for MuSiQue. Tab. 3 reports mean exact match and F1 score on the provided 1000 samples for each dataset. These results show π-CoT outperforming both Standard RAG and HippoRAG 2, demonstrating that offline fact extraction is not necessary for good accuracy on multi-hop question-answering tasks.

## 5.2 IN-CONTEXT QUESTION-ANSWERING

When all the necessary information fits in the context window of the model, a natural question is *when does formal reasoning (via π-CoT) benefit reasoning in natural-language (via CoT)?*. To investigate this, we employ the distractor variants of HotpotQA, 2WikiMultiHopQA, and MuSiQue, where the gold passages are presented in-context alongside a small number of irrelevant (“distractor”) passages. We also include two versions of the PhantomWiki benchmark from Gong et al. (2025): PW-S and PW-M. Unlike the other three datasets that are curated from Wikipedia, PhantomWiki generates challenging multi-hop questions from fictional universes to ensure contamination-free LLM evaluation. We provide supplementary experiment details in Sec. C.2.

Tab. 4(a) shows that on HotpotQA, 2WikiMultiHopQA, and MuSiQue, there is no significant difference in accuracy between CoT and π-CoT. Providing the gold passages to model makes the task considerably easier than considering all of Wikipedia (Min et al., 2019). Thus, Llama-3.3-70B-Instruct

<sup>7</sup>Fine-tuned approaches are currently the state-of-the-art on HotpotQA, 2WikiMultiHopQA, and MuSiQue.

378  
 379  
 380  
 381  
 382  
 383  
 384  
 385  
 386  
 387  
 388  
 389  
 390  
 391  
 392  
 393  
 394  
 395  
 396  
 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  
 429  
 430  
 431  
 Table 4: **Accuracy of prompting-based methods on in-context QA.** We report exact match (EM) and F1 score on HotpotQA, 2WikiMultiHopQA, MuSiQue, PhantomWiki with a corpus size of 50 articles (PW-S), and PhantomWiki with a corpus size of 500 articles (PW-M). Bold indicates that  $\pi$ -CoT significantly outperforms CoT ( $p < 0.05$ ) according to a paired samples  $t$ -test given a dataset and metric. We report the computational cost in Tab. 8.

Method	HotpotQA		2WikiMultiHopQA		MuSiQue	
	EM $\uparrow$	F1 $\uparrow$	EM $\uparrow$	F1 $\uparrow$	EM $\uparrow$	F1 $\uparrow$
<i>Llama-3.3-70B-Instruct</i>						
CoT	$62.4 \pm 2.2$	$77.4 \pm 1.6$	$76.4 \pm 1.9$	$83.9 \pm 1.5$	$51.2 \pm 2.2$	$63.7 \pm 1.9$
$\pi$ -CoT	$58.0 \pm 2.2$	$77.7 \pm 1.5$	$72.8 \pm 2.0$	$82.5 \pm 1.5$	$46.4 \pm 2.2$	$63.1 \pm 1.9$
<i>DeepSeek-R1-Distill-Qwen-32B</i>						
CoT	$57.0 \pm 2.2$	$74.2 \pm 1.7$	$75.2 \pm 1.9$	$82.9 \pm 1.6$	$45.8 \pm 2.2$	$57.3 \pm 2.0$
$\pi$ -CoT	$56.8 \pm 2.2$	$75.2 \pm 1.6$	$74.6 \pm 1.9$	$84.2 \pm 1.5$	$46.0 \pm 2.2$	$59.6 \pm 2.0$
(a) Real-world (Wikipedia-based) multi-hop QA datasets						
Method	PW-S		PW-M			
	EM $\uparrow$	F1 $\uparrow$	EM $\uparrow$	F1 $\uparrow$		
<i>Llama-3.3-70B-Instruct</i>						
CoT	$52.8 \pm 1.3$	$71.9 \pm 1.0$	$27.5 \pm 1.2$	$41.7 \pm 1.1$		
$\pi$ -CoT	<b><math>78.8 \pm 1.1</math></b>	<b><math>91.4 \pm 0.6</math></b>	<b><math>31.1 \pm 1.2</math></b>	<b><math>56.9 \pm 1.0</math></b>		
<i>DeepSeek-R1-Distill-Qwen-32B</i>						
CoT	$54.4 \pm 1.3$	$75.6 \pm 0.9$	$17.5 \pm 1.0$	$28.0 \pm 1.0$		
$\pi$ -CoT (Ours)	<b><math>82.7 \pm 1.0</math></b>	<b><math>88.2 \pm 0.8</math></b>	$18.4 \pm 1.0$	<b><math>31.1 \pm 1.0</math></b>		
(b) Synthetic multi-hop QA datasets						

with CoT may be hitting a performance ceiling. On PW-S, Llama-3.3-70B-Instruct with CoT achieves an F1 score of  $71.9 \pm 1.0\%$ . This increases to  $91.4 \pm 0.6\%$  with  $\pi$ -CoT. On PW-M, which has a corpus 10 times the size of PW-S, the performance of Llama-3.3-70B-Instruct with CoT drops to  $41.7 \pm 1.1\%$  F1. We posit this drop is mainly due to the inherent challenges of long-context retrieval. On PW-M,  $\pi$ -CoT significantly boosts the F1 score to  $56.9 \pm 1.0$  F1—a relative gain of 36%. In Tab. 4(b), we report similar findings using Deepseek-R1-Distill-Qwen-32B as the language model.

Finally, to understand where the gains on PW-S and PW-M come from, we plot accuracy versus the ground-truth difficulty level that comes associated with each question. [Gong et al. \(2025\)](#) defines this difficulty level as the number of hops required to answer the question. According to Fig. 3, the accuracy of  $\pi$ -CoT and CoT is similar for questions with low difficulty. However,  $\pi$ -CoT diverges from CoT as the difficulty increases. Since higher difficulty questions require traversing more reasoning paths than lower difficulty questions, our results show that  $\pi$ -CoT is better able to keep track of multi-hop, multi-branch reasoning than CoT.

## 6 ABLATION ANALYSIS

In this section, we want to analyze how each component in  $\pi$ -CoT contributes to answering the multi-hop question.

**Contributions of the Prolog component.** Note that when Prolog execution succeeds, the LLM in the final CoT step is instructed to simply copy the Prolog answer as the final output. This setup naturally induces the following categorization: Prolog execution either returns a non-empty solution set ( $S_T \neq \emptyset$ ) or an empty one ( $S_T = \emptyset$ ), and the final CoT answer is either correct ( $\pi$ -CoT ✓) or

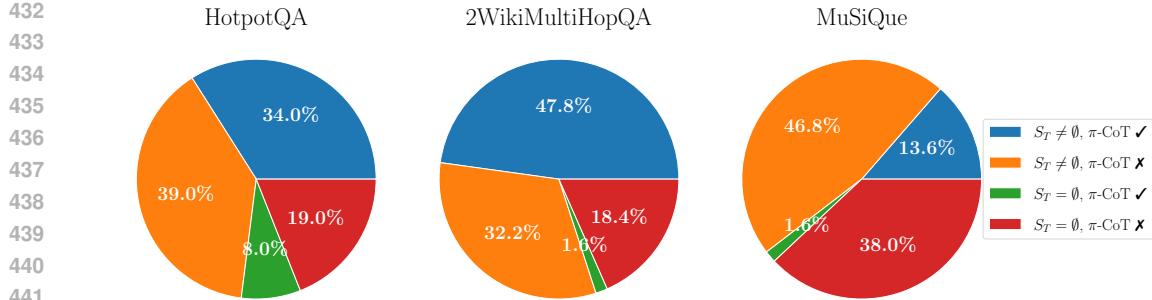


Figure 4: **Analysis of Prolog component based on execution success.**  $S_T \neq \emptyset$  means that Prolog outputs an answer, while  $S_T = \emptyset$  means that Prolog does not output an answer.  $\pi\text{-CoT} \checkmark$  means that  $\pi\text{-CoT}$  generated the ground-truth answer exactly, while  $\pi\text{-CoT} \times$  means that  $\pi\text{-CoT}$  generated an incorrect answer. We use the results from Tab. 1.

incorrect ( $\pi\text{-CoT} \times$ ). Here, an answer is considered correct only if it exactly matches the ground-truth answer.

Fig. 4 shows a breakdown of the predictions from Tab. 1 across these four cases. Across datasets, Prolog execution returns an answer in 73% (HotpotQA), 80% (2WikiMultiHopQA), and 60.4% (MuSiQue) of cases. When Prolog returns an answer,  $\pi\text{-CoT}$  is consistently more accurate (46.6%, 60.0%, 22.5%) than the best baseline method in Tab. 1 (40.4%, 37.2%, 17.6%). Even when Prolog sometimes does not return an answer, using the artifacts coming from the Prolog execution can help lead to a correct final answer. Notably, 7.8% of the questions in HotpotQA can be answered correctly this way. For the remaining (27%, 20%, 39.6%) cases, we provide a detailed categorization of the failures in App. E.

**Contributions of components in the final CoT prompt.** We remove specific components from the prompt when generating the final answer (Sec. 4). Tab. 5 shows that removing the passages leads to a (4.7%, 1.2%, 2.4%) drop in F1. This suggests that most of the information needed to answer the question already lies in the notes and Prolog answer. We next remove the notes and see performance drop by (4.8%, 0.6%, 0.8%) F1. On HotpotQA especially, notes are important for fuzzy matching when exact string matching fails. Finally, we remove the Prolog answer and see performance drop to near zero for 2WikiMultiHopQA and MuSiQue. This demonstrates the utility of the Prolog answer to the final chain-of-thought reasoning step. Interestingly, Llama-3.3-70B-Instruct is able to achieve 24.1% F1 solely using its internal knowledge. We did not find that including the passages to pose a challenge with long context. Thus, we included all three components for our final model for the best accuracy.

## 7 CONCLUSIONS & FUTURE WORK

In this work, we investigate how formal reasoning can guide reasoning in natural language. We introduce  $\pi\text{-CoT}$ , a novel prompting strategy that initializes the context of an LLM with the intermediate outputs of Prolog-guided execution. Our results show that even strong LLMs like

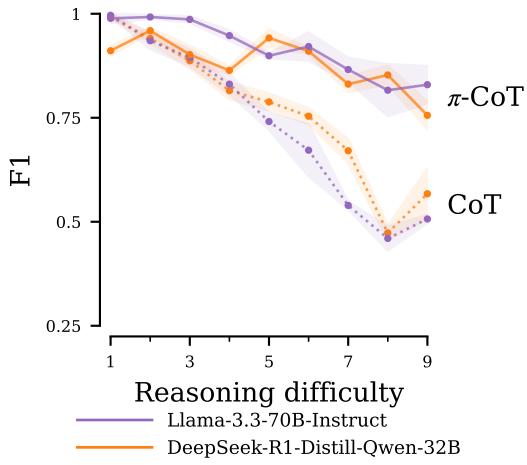


Figure 3: **F1 score vs. difficulty, as measured by number of reasoning steps.** We use the synthetic PW-S benchmark from Gong et al. (2025) and display mean  $\pm 1$  standard error. For each model, we evaluate CoT and  $\pi\text{-CoT}$  prompting.

On HotpotQA especially, notes are important for fuzzy matching when exact string matching fails. Finally, we remove the Prolog answer and see performance drop to near zero for 2WikiMultiHopQA and MuSiQue. This demonstrates the utility of the Prolog answer to the final chain-of-thought reasoning step. Interestingly, Llama-3.3-70B-Instruct is able to achieve 24.1% F1 solely using its internal knowledge. We did not find that including the passages to pose a challenge with long context. Thus, we included all three components for our final model for the best accuracy.

486  
487 Table 5: **Ablation analysis of the components in the final chain-of-thought reasoning step.** We  
488 use the experimental setup from Tab. 1.  
489

Method	HotpotQA		2WikiMultiHopQA		MuSiQue	
	EM $\uparrow$	F1 $\uparrow$	EM $\uparrow$	F1 $\uparrow$	EM $\uparrow$	F1 $\uparrow$
$\pi$ -CoT	42.0	59.1	49.4	57.5	15.2	25.7
w/o Passages	37.2	54.4	48.4	56.3	12.6	23.3
w/o Notes	34.0	49.6	47.8	55.7	12.4	22.5
w/o Prolog Answer	19.2	24.1	1.0	1.0	3.2	4.2

497  
498 Llama-3.3-70B-Instruct and Deepseek-R1-Distill-Qwen-32B benefit from being guided through  
499 structured reasoning steps, especially on complex, multi-step tasks. More broadly, our work  
500 demonstrates the potential of bringing explicit planning and state tracking into language model  
501 behavior.

502  
503 Despite being a purely prompting-based strategy,  $\pi$ -CoT has potential implications for training future  
504 language models to serve as agents that can piece together knowledge across large, dynamic corpora.  
505 Inspired by DeepSeek R1 (Guo et al., 2025), significant effort has been made to couple reasoning  
506 with retrieval using reinforcement learning (Li et al., 2025; Jin et al., 2025a; Song et al., 2025). An  
507 interesting future direction is training language models to generate structured queries instead.  $\pi$ -CoT  
508 provides a way to execute these queries without a pre-existing database. Some steps in  $\pi$ -CoT may  
509 not even require calling an LLM, if the information resides directly in a pre-existing database. Thus,  
510 enabling  $\pi$ -CoT to leverage both unstructured and structured data is a natural next step.

## 511 ETHICS STATEMENT

512 Our work adheres to the ICLR Code of Ethics, and does not pose any societal, personal, or  
513 organizational risks.

## 514 REPRODUCIBILITY STATEMENT

515 To encourage reproducibility, we use free and open-source software and LLMs. We also include  
516 details of our experimental setups in Sec. C. We report sample sizes, standard errors, and random  
517 seeds where possible.

## 518 REFERENCES

519 Gabor Angeli, Melvin Jose Johnson Premkumar, and Christopher D Manning. Leveraging linguistic  
520 structure for open domain information extraction. In *Proceedings of the 53rd Annual Meeting*  
521 *of the Association for Computational Linguistics and the 7th International Joint Conference on*  
522 *Natural Language Processing (Volume 1: Long Papers)*, pages 344–354, 2015. (Cited on page 2.)

523 Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning to  
524 retrieve, generate, and critique through self-reflection. In *The Twelfth International Conference on*  
525 *Learning Representations*, 2023. (Cited on pages 1 and 6.)

526 Nasim Borazjanizadeh and Steven T Piantadosi. Reliable reasoning beyond natural language. *arXiv*  
527 *preprint arXiv:2407.11373*, 2024. (Cited on page 2.)

528 Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. The snli corpus.  
529 2015. (Cited on page 15.)

530 Jifan Chen and Greg Durrett. Understanding dataset design choices for multi-hop reasoning. *arXiv*  
531 *preprint arXiv:1904.12106*, 2019. (Cited on page 7.)

540 Xinyun Chen, Chen Liang, Adams Wei Yu, Denny Zhou, Dawn Song, and Quoc V Le. Neural  
 541 symbolic reader: Scalable integration of distributed and symbolic representations for reading  
 542 comprehension. In *International Conference on Learning Representations*, 2019. (Cited on page 2.)  
 543

544 William F Clocksin and Christopher S Mellish. *Programming in PROLOG*. Springer Science &  
 545 Business Media, 2003. (Cited on page 1.)

546 Alain Colmerauer and Philippe Roussel. The birth of prolog. In *History of programming  
 547 languages—II*, pages 331–367. 1996. (Cited on page 3.)

548 Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine Li, Liwei Jiang, Bill Yuchen Lin, Sean  
 549 Welleck, Peter West, Chandra Bhagavatula, Ronan Le Bras, et al. Faith and fate: Limits of  
 550 transformers on compositionality. *Advances in Neural Information Processing Systems*, 36:  
 551 70293–70332, 2023. (Cited on pages 1 and 2.)

552 Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt,  
 553 Dasha Metropolitansky, Robert Osazuwa Ness, and Jonathan Larson. From local to global: A  
 554 graph rag approach to query-focused summarization. *arXiv preprint arXiv:2404.16130*, 2024.  
 555 (Cited on pages 3 and 7.)

556 Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yixin Dai, Jiawei Sun,  
 557 Haofen Wang, and Haofen Wang. Retrieval-augmented generation for large language models: A  
 558 survey. *arXiv preprint arXiv:2312.10997*, 2(1), 2023. (Cited on page 2.)

559 Albert Gong, Kamilė Stankevičiūtė, Chao Wan, Anmol Kabra, Raphael Thesmar, Johann Lee,  
 560 Julius Klenke, Carla P Gomes, and Kilian Q Weinberger. Phantomwiki: On-demand datasets for  
 561 reasoning and retrieval evaluation. In *Forty-second International Conference on Machine Learning*,  
 562 2025. (Cited on pages 7, 8, and 9.)

563 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad  
 564 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of  
 565 models. *arXiv preprint arXiv:2407.21783*, 2024. (Cited on page 6.)

566 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu,  
 567 Ruoyu Zhang, Shirong Ma, Xiao Bi, et al. Deepseek-r1 incentivizes reasoning in llms through  
 568 reinforcement learning. *Nature*, 645(8081):633–638, 2025. (Cited on pages 6 and 10.)

569 Zirui Guo, Lianghao Xia, Yanhua Yu, Tu Ao, and Chao Huang. Lightrag: Simple and fast  
 570 retrieval-augmented generation. *arXiv preprint arXiv:2410.05779*, 2024. (Cited on pages 3 and 7.)

571 Bernal Jiménez Gutiérrez, Yiheng Shu, Weijian Qi, Sizhe Zhou, and Yu Su. From rag to memory:  
 572 Non-parametric continual learning for large language models. In *Forty-second International  
 573 Conference on Machine Learning*, 2025. (Cited on pages 3 and 7.)

574 Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. Constructing a multi-hop  
 575 qa dataset for comprehensive evaluation of reasoning steps. In *Proceedings of the 28th International  
 576 Conference on Computational Linguistics*, pages 6609–6625, 2020. (Cited on page 6.)

577 Bernal Jimenez Gutierrez, Yiheng Shu, Yu Gu, Michihiro Yasunaga, and Yu Su. Hipporag:  
 578 Neurobiologically inspired long-term memory for large language models. *Advances in Neural  
 579 Information Processing Systems*, 37:59532–59569, 2024. (Cited on page 3.)

580 Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Zamani, and  
 581 Jiawei Han. Search-r1: Training llms to reason and leverage search engines with reinforcement  
 582 learning. *arXiv preprint arXiv:2503.09516*, 2025a. (Cited on pages 2 and 10.)

583 Jiajie Jin, Yutao Zhu, Zhicheng Dou, Guanting Dong, Xinyu Yang, Chenghao Zhang, Tong Zhao,  
 584 Zhao Yang, and Ji-Rong Wen. Flashrag: A modular toolkit for efficient retrieval-augmented  
 585 generation research. In *Companion Proceedings of the ACM on Web Conference 2025*, pages  
 586 737–740, 2025b. (Cited on page 6.)

587 Omar Khattab, Keshav Santhanam, Xiang Lisa Li, David Hall, Percy Liang, Christopher Potts,  
 588 and Matei Zaharia. Demonstrate-search-predict: Composing retrieval and language models for  
 589 knowledge-intensive nlp. *arXiv preprint arXiv:2212.14024*, 2022. (Cited on page 2.)

594 Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish  
 595 Sabharwal. Decomposed prompting: A modular approach for solving complex tasks. In *The*  
 596 *Eleventh International Conference on Learning Representations*, 2023. (Cited on pages 1 and 2.)  
 597

598 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large  
 599 language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:  
 600 22199–22213, 2022. (Cited on page 1.)

601 Robert Kowalski and Steve Smoliar. Logic for problem solving. *ACM SIGSOFT Software Engineering*  
 602 *Notes*, 7(2):61–62, 1982. (Cited on page 1.)

603 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph  
 604 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model  
 605 serving with pagedattention. In *Proceedings of the 29th symposium on operating systems principles*,  
 606 pages 611–626, 2023. (Cited on page 17.)

607

608 Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan Raiman, Mohammad Shoeybi, Bryan Catanzaro,  
 609 and Wei Ping. Nv-embed: Improved techniques for training llms as generalist embedding models.  
 610 In *The Thirteenth International Conference on Learning Representations*, 2024. (Cited on page 7.)

611 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal,  
 612 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented  
 613 generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*,  
 614 33:9459–9474, 2020. (Cited on page 2.)

615

616 Xiaoxi Li, Guanting Dong, Jiajie Jin, Yuyao Zhang, Yujia Zhou, Yutao Zhu, Peitian Zhang, and  
 617 Zhicheng Dou. Search-o1: Agentic search-enhanced large reasoning models. *arXiv preprint*  
 618 *arXiv:2501.05366*, 2025. (Cited on pages 2 and 10.)

619 Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and  
 620 Percy Liang. Lost in the middle: How language models use long contexts. *Transactions of the*  
 621 *Association for Computational Linguistics*, 12:157–173, 2024. (Cited on page 1.)

622 Robert Lo, Abishek Sridhar, Frank F Xu, Hao Zhu, and Shuyan Zhou. Hierarchical prompting  
 623 assists large language model on web navigation. In *Findings of the Association for Computational*  
 624 *Linguistics: EMNLP 2023*, pages 10217–10244, 2023. (Cited on page 1.)

625

626 John McCarthy et al. *Programs with common sense*. RLE and MIT computation center Cambridge,  
 627 MA, USA, 1960. (Cited on page 1.)

628

629 Sewon Min, Eric Wallace, Sameer Singh, Matt Gardner, Hannaneh Hajishirzi, and Luke Zettlemoyer.  
 630 Compositional questions do not necessitate multi-hop reasoning. *arXiv preprint arXiv:1906.02900*,  
 631 2019. (Cited on pages 1 and 7.)

632

633 Pruthvi Patel, Swaroop Mishra, Mihir Parmar, and Chitta Baral. Is a question decomposition unit  
 634 all we need? In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language*  
 635 *Processing*, pages 4553–4569, 2022. (Cited on page 1.)

636

637 Kevin Pei, Ishan Jindal, Kevin Chen-Chuan Chang, ChengXiang Zhai, and Yunyao Li. When to use  
 638 what: An in-depth comparative empirical analysis of openie systems for downstream applications.  
 639 In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*  
 640 (*Volume 1: Long Papers*), pages 929–949, 2023. (Cited on page 3.)

641

642 Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A Smith, and Mike Lewis. Measuring  
 643 and narrowing the compositionality gap in language models. In *Findings of the Association for*  
 644 *Computational Linguistics: EMNLP 2023*, pages 5687–5711, 2023. (Cited on pages 1 and 2.)

645

646 Stephen Robertson, Hugo Zaragoza, et al. The probabilistic relevance framework: Bm25 and beyond.  
 647 *Foundations and Trends® in Information Retrieval*, 3(4):333–389, 2009. (Cited on page 6.)

648

649 John Alan Robinson. A machine-oriented logic based on the resolution principle. *Journal of the*  
 650 *ACM (JACM)*, 12(1):23–41, 1965. (Cited on page 1.)

651

652 JK Rowling. *Harry Potter*, volume 2007. London Bloomsbury, 1997. (Cited on page 3.)

648 Stuart Russell, Peter Norvig, and Artificial Intelligence. A modern approach. *Artificial Intelligence.*  
 649 *Prentice-Hall, Englewood Cliffs*, 25(27):79–80, 1995. (Cited on pages 1 and 3.)  
 650

651 Parth Sarthi, Salman Abdulla, Aditi Tuli, Shubh Khanna, Anna Goldie, and Christopher D Manning.  
 652 Raptor: Recursive abstractive processing for tree-organized retrieval. In *The Twelfth International*  
 653 *Conference on Learning Representations*, 2024. (Cited on pages 3 and 7.)

654 Herbert A Simon and Allen Newell. Human problem solving: The state of the theory in 1970.  
 655 *American psychologist*, 26(2):145, 1971. (Cited on page 1.)  
 656

657 Huatong Song, Jinhao Jiang, Yingqian Min, Jie Chen, Zhipeng Chen, Wayne Xin Zhao, Lei Fang, and  
 658 Ji-Rong Wen. R1-searcher: Incentivizing the search capability in llms via reinforcement learning.  
 659 *arXiv preprint arXiv:2503.05592*, 2025. (Cited on pages 2 and 10.)

660 Leon Sterling and Ehud Y Shapiro. *The art of Prolog: advanced programming techniques*. MIT  
 661 press, 1994. (Cited on page 3.)  
 662

663 Dídac Surís, Sachit Menon, and Carl Vondrick. Vipergrpt: Visual inference via python execution for  
 664 reasoning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages  
 665 11888–11898, 2023. (Cited on page 2.)

666 Xiaoyu Tan, Yongxin Deng, Xihe Qiu, Weidi Xu, Chao Qu, Wei Chu, Yinghui Xu, and Yuan Qi.  
 667 Thought-like-pro: Enhancing reasoning of large language models through self-driven prolog-based  
 668 chain-of-thought. *arXiv preprint arXiv:2407.14562*, 2024. (Cited on page 2.)  
 669

670 Hieu Tran, Zonghai Yao, Junda Wang, Yifan Zhang, Zhichao Yang, and Hong Yu. Rare:  
 671 Retrieval-augmented reasoning enhancement for large language models. *arXiv preprint*  
 672 *arXiv:2412.02830*, 2024. (Cited on page 2.)

673 Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. MuSiQue: Multihop  
 674 questions via single-hop question composition. *Transactions of the Association for Computational*  
 675 *Linguistics*, 10:539–554, 2022. (Cited on page 6.)  
 676

677 Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. Interleaving retrieval  
 678 with chain-of-thought reasoning for knowledge-intensive multi-step questions. In *Proceedings*  
 679 *of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long*  
 680 *Papers)*, pages 10014–10037, 2023. (Cited on pages 2, 7, and 17.)

681 Priyesh Vakharia, Abigail Kufeldt, Max Meyers, Ian Lane, and Leilani H Gilpin. Proslm: A  
 682 prolog synergized language model for explainable domain specific knowledge based question  
 683 answering. In *International Conference on Neural-Symbolic Learning and Reasoning*, pages  
 684 291–304. Springer, 2024. (Cited on page 2.)  
 685

686 Liang Wang, Haonan Chen, Nan Yang, Xiaolong Huang, Zhicheng Dou, and Furu Wei.  
 687 Chain-of-retrieval augmented generation. *ArXiv*, abs/2501.14342, 2025. URL <https://api.semanticscholar.org/CorpusID:275906944>. (Cited on page 2.)  
 688

689 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha  
 690 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language  
 691 models. *arXiv preprint arXiv:2203.11171*, 2022. (Cited on page 1.)  
 692

693 Leon Weber, Pasquale Minervini, Jannes Münchmeyer, Ulf Leser, and Tim Rocktäschel. Nlprolog:  
 694 Reasoning with weak unification for question answering in natural language. In *Proceedings*  
 695 *of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for  
 696 Computational Linguistics, 2019. (Cited on page 2.)

697 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny  
 698 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*  
 699 *neural information processing systems*, 35:24824–24837, 2022. (Cited on pages 1 and 2.)  
 700

701 Jan Wielemaker, Tom Schrijvers, Markus Triska, and Torbjörn Lager. Swi-prolog. *Theory and*  
 702 *Practice of Logic Programming*, 12(1-2):67–96, 2012. (Cited on page 3.)

702 Tomer Wolfson, Mor Geva, Ankit Gupta, Matt Gardner, Yoav Goldberg, Daniel Deutch, and Jonathan  
 703 Berant. Break it down: A question understanding benchmark. *Transactions of the Association for  
 704 Computational Linguistics*, 8:183–198, 2020. (Cited on page 1.)  
 705

706 Katherine Wu and Yanhong A Liu. Lp-Im: No hallucinations in question answering with logic  
 707 programming. *arXiv preprint arXiv:2502.09212*, 2025. (Cited on page 2.)  
 708

709 Xiaocheng Yang, Bingsen Chen, and Yik-Cheung Tam. Arithmetic reasoning with llm: Prolog  
 710 generation & permutation. In *Proceedings of the 2024 Conference of the North American Chapter  
 711 of the Association for Computational Linguistics: Human Language Technologies (Volume 2:  
 712 Short Papers)*, pages 699–710, 2024. (Cited on page 2.)  
 713

714 Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov,  
 715 and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question  
 716 answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language  
 717 Processing*, pages 2369–2380, 2018. (Cited on page 6.)  
 718

719 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan  
 720 Cao. React: Synergizing reasoning and acting in language models. In *The Eleventh International  
 721 Conference on Learning Representations*, 2023. (Cited on pages 1 and 2.)  
 722

723 Howard Yen, Tianyu Gao, Minmin Hou, Ke Ding, Daniel Fleischer, Peter Izsak, Moshe Wasserblat,  
 724 and Danqi Chen. Helmet: How to evaluate long-context language models effectively and thoroughly.  
 725 *arXiv preprint arXiv:2410.02694*, 2024. (Cited on page 1.)  
 726

727 John M Zelle and Raymond J Mooney. Learning to parse database queries using inductive  
 728 logic programming. In *Proceedings of the national conference on artificial intelligence*, pages  
 729 1050–1055, 1996. (Cited on page 3.)  
 730

731 Luke S Zettlemoyer and Michael Collins. Learning to map sentences to logical form: Structured  
 732 classification with probabilistic categorial grammars. *arXiv preprint arXiv:1207.1420*, 2012. (Cited  
 733 on page 3.)  
 734

735 Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in  
 736 large language models. *arXiv preprint arXiv:2210.03493*, 2022. (Cited on page 1.)  
 737

738 Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans,  
 739 Claire Cui, Olivier Bousquet, Quoc V Le, et al. Least-to-most prompting enables complex  
 740 reasoning in large language models. In *The Eleventh International Conference on Learning  
 741 Representations*, 2023. (Cited on pages 1 and 2.)  
 742

743 Shaowen Zhou, Bowen Yu, Aixin Sun, Cheng Long, Jingyang Li, Haiyang Yu, Jian Sun, and Yongbin  
 744 Li. A survey on neural open information extraction: Current status and future directions. *arXiv  
 745 preprint arXiv:2205.11725*, 2022. (Cited on page 3.)  
 746

747

748

749

750

751

752

753

754

755

756 A IMPLEMENTATION DETAILS OF SLICE MODULE  
757758 **Inputs.** At step  $t \in \{1, 2, \dots, T\}$ , SLICE<sup>8</sup> takes as input:  
759760 

- 761 • the sub-query,  $q_t$ ;
- 762 • the solution from the previous step,  $S_{t-1}$ ; and
- 763 • the document corpus,  $\mathcal{C}$ .

764 At step  $t = 1$ , there are no solutions, so  $S_0 = \emptyset$ . At step  $t > 1$ , we assign any values in  $S_{t-1}$  to the  
765 variables in  $q_t$ . Let's take the second sub-query,  $q_2$ , from Fig. 1 as an example. As shown in Fig. 2,  
766 the previous solution  $S_1$  assigns the following values to X: "Quirrell," "Lockhart," "Lupin," etc. These  
767 must be substituted into  $q_2 = \text{werewolf}(X)$ , yielding the queries  $\text{werewolf}(\text{"Quirrell"})$ ,  
768  $\text{werewolf}(\text{"Lockhart"})$ ,  $\text{werewolf}(\text{"Lupin"})$ , etc. to resolve for the current step.  
769770 While the query  $\text{werewolf}(\text{"Quirrell"})$  provides a succinct way of checking whether the  
771 claim, "Quirrell was a werewolf," is true, it's still uninterpretable to an LLM. Thus, we need a  
772 mechanism to translate Prolog facts to natural-language statements. Our solution is to prompt<sup>9</sup> the  
773 LLM to generate a **definition** for  $q_t$ . Each definition comprises two templates:  
774775 

- 776 • A **question template**, which maps  $q_t$  to a question (e.g., "Who are the Defense against the Dark  
777 Arts teachers at Hogwarts?").
- 778 • A **statement template**, which maps  $q_t$  to a claim (e.g., "Lockhart is a Defense against the Dark  
779 Arts teacher at Hogwarts").

780 Each template serves a different purpose, either *fact extraction* or *fact verification*.  
781782 **Fact extraction.** When  $q_t$  introduces a new unassigned variable, SLICE uses the question template  
783 to map  $q_t$  to a natural-language question (see **Prolog**→**NL** in Fig. 2). Next, we call the LLM using  
784 chain-of-thought prompting to answer this question (see **LLM** in Fig. 2). We provide the LLM  
785 prompt in Sec. B.2. In this work, we consider two strategies to retrieve the relevant evidence for the  
786 question:  
787788 (S1) In the **RAG setting**, we use an external retriever to obtain the top- $k$  passages from  $\mathcal{C}$ .  
789790 (S2) In the **in-context setting**, all passages of  $\mathcal{C}$  are provided in the prompt.  
791792 Using its innate reading comprehension abilities, the LLM locates the answer to the question from the  
793 provided passages and responds with (potentially) multiple answers (e.g., Quirrell, Lockhart, Lupin,  
794 etc). These answers are parsed back into Prolog facts (see **Prolog**←**NL** in Fig. 2) and added to the  
795 knowledge base.  
796797 **Fact verification.** When all variables in  $q_t$  can be assigned values from  $S_{t-1}$ , we must check that the  
798 fact is true. SLICE uses the statement template to generate an entailment question (Bowman et al.,  
799 2015). For example,  $\text{werewolf}(\text{"Quirrell"})$  corresponds to the question, "Is the following  
800 statement true or false? Quirrell was a werewolf." Equipped with S1 or S2, the LLM answers this  
801 question (see **LLM** in Fig. 2). A true claim is added to the knowledge base as a fact; a false claim  
802 is simply ignored (and nothing is added to knowledge base). In our running example, the only fact  
803 added to the knowledge base at step  $t = 2$  is  $\text{werewolf}(\text{"Lupin"})$  (see green text in Fig. 2).  
804805 **Output.** Finally, the output of SLICE is the updated solution  $S_t$ , which can be obtained by querying  
806 the knowledge base with  $(q_1, \dots, q_t)$ .  
807808 B PROMPT TEMPLATES  
809810 B.1 PROLOG QUERY AND DEFINITIONS GENERATION  
811812 We use the following prompt template for the Prolog query generation of Sec. 4:  
8138<sup>8</sup>short for Single-step Logical Inference with Contextual Evidence9<sup>9</sup>We provide the definitions generation prompt in Sec. B.1.

810 You will be provided a question. Your goal is to devise a  
 811 Prolog query to answer this question. Your response must end in  
 812 "`**Query:** <query>\n**Target:** <target>\n**Definition:**`  
 813 `<definition>`", where `<query>` is a Prolog query that when  
 814 executed, will yield the answer to the question, `<target>`  
 815 is the target variable in the Prolog query to be returned  
 816 as the final answer, and `<definition>` defines the semantic  
 817 meaning of predicates in the Prolog query.  
 818  
 819 Here are some examples:  
 820 (START OF EXAMPLES)  
 821 {examples}  
 822 (END OF EXAMPLES)  
 823  
 824 Question: {question}  
 825 Answer:  
 826  
 827 To form the prompt, `examples` is replaced with few-shot examples specific to each dataset and  
 828 question is replaced with the natural-language question.  
 829  
 830 **B.2 CHAIN-OF-THOUGHT FACT EXTRACTION AND VERIFICATION**  
 831  
 832 You are given the following evidence:  
 833 (BEGIN EVIDENCE)  
 834 {{evidence}}  
 835 (END EVIDENCE)  
 836  
 837 You will be provided a question. If there is a single answer,  
 838 your response must end with the final answer enclosed in tags:  
 839 `<answer>FINAL_ANSWER</answer>`  
 840 If there are multiple answers, your response must end with the  
 841 final answers enclosed in tags:  
 842 `<answer>FINAL_ANSWER_1, FINAL_ANSWER_2, ..., FINAL_ANSWER_N</answer>`.  
 843 If `FINAL_ANSWER_N` is a string, it must be enclosed in double quotes.  
 844 For example, `<answer>"FINAL_ANSWER_1", "FINAL_ANSWER_2"</answer>`  
 845 If `FINAL_ANSWER_N` is a date, it must be formatted as  
 846 `date(year, month, day)`.  
 847 If no information is available to answer the question,  
 848 your response must end with: `<answer></answer>`.  
 849  
 850 Here are some examples:  
 851 (START OF EXAMPLES)  
 852 {{examples}}  
 853 (END OF EXAMPLES)  
 854  
 855 Question: {{question}}  
 856 Answer:  
 857  
 858 To instantiate the prompt, `evidence` is replaced by relevant passages, `examples` is replaced with  
 859 dataset-specific few-shot examples, and `question` is replaced with the natural-language question  
 860 (in the case of fact extraction) or entailment question (in the case of fact verification). Both the  
 861 instructions and few-shot examples encourage the model to format the answer as Prolog literals so  
 862 that they can be properly inserted into the Prolog knowledge base.  
 863  
 864 **B.3  $\pi$ -CoT**  
 865  
 866 You are given the following information:  
 867 (BEGIN NOTES)  
 868 {{notes}}

```

864 (END NOTES)
865
866 (BEGIN EVIDENCE)
867 {{evidence}}
868 (END EVIDENCE)
869
870 You will be provided a question and an answer from a previous
871 attempt. If the previous answer is not empty
872 (e.g. <answer>...</answer>), you should copy the answer directly.
873 If the previous answer is empty (i.e. <answer></answer>),
874 you should try to answer the question using
875 the notes and evidence provided. If there is a single answer,
876 your response must end with the final answer enclosed in tags:
877 <answer>FINAL_ANSWER</answer>
878 If there are multiple answers, your response must end with the
879 final answers enclosed in tags:
880 <answer>FINAL_ANSWER_1,FINAL_ANSWER_2,...,FINAL_ANSWER_N</answer>.
881 If no information is available to answer the question,
882 your response must end with: <answer></answer>.

883 Here are some examples:
884 (START OF EXAMPLES)
885 {{examples}}
886 (END OF EXAMPLES)
887 Question: {{question}}
888 Previous Answer: <answer>{{answer}}</answer>
889 Answer:
890
891

```

## C SUPPLEMENTARY EXPERIMENT DETAILS

### C.1 FULLWIKI EXPERIMENT DETAILS

**Retrieval setup.** We use the `wiki18_100w` corpus from [https://huggingface.co/datasets/RUC-NLPIR/FlashRAG\\_datasets](https://huggingface.co/datasets/RUC-NLPIR/FlashRAG_datasets) and use the code from <https://github.com/RUC-NLPIR/FlashRAG> to build our BM25 index. We allow  $k = 14$ ,  $k = 16$ , and  $k = 8$  chunks per retrieval call for HotpotQA, 2WikiMultiHopQA, and MuSiQue, respectively.

**Baseline implementations.** For standard RAG, we use the Python implementation of CoTRAGAgent from <https://github.com/kilian-group/phantom-wiki> and write few-shot examples for each dataset. For Self-Ask, we use the Python implementation and few-shot examples from <https://github.com/RUC-NLPIR/FlashRAG>. For IRCoT, we use the Python implementation from <https://github.com/RUC-NLPIR/FlashRAG> and the GPT3 (code-davinci-002) few-shot examples from <https://github.com/StonyBrookNLP/ircot> (see also (Trivedi et al., 2023, App. G)). We set the maximum iterations for Self-Ask and IRCOT to be 4.

**LLM configuration.** We run Llama-3.3-70B-Instruct on 8 A6000s using vLLM (Kwon et al., 2023) and use greedy decoding with maximum generation tokens 4096. We use the full 128K context length.

### C.2 DISTRACTOR EXPERIMENT DETAILS

**LLM configuration.** For the Llama-3.3-70B-Instruct results, we use vLLM running on 8 A6000s and use greedy decoding with maximum generation tokens 4096. For the Deepseek-R1-Distill-Qwen-32B results, we use vLLM running on 6 A6000s and use sampling temperature 0.6, top-p 0.95, and max generation tokens 16384. We use the full 128K context length for both models.

**PhantomWiki dataset.** For the experiment of Tab. 4(b), we use the code at <https://github.com/kilian-group/phantom-wiki> to generate two synthetic multi-hop QA datasets. Tab. 6 lists the configurations for PW-S and PW-M. Each dataset has 1500 questions.

Table 6: Configurations for PhantomWiki dataset generation.

Parameter	PW-S	PW-M
Question depth	20	20
Number of family trees	10	100
Max family tree size	50	50
Max family tree depth	20	20
Mode	Easy	Easy
Number of questions per template	10	10
Seeds	{1,2,3}	{1,2,3}

### C.3 LLM USAGE STATEMENT

Large language models were used for proofreading, revising, and literature search. All claims and arguments were drafted and verified by the authors.

## D ADDITIONAL DETAILS OF COMPUTATIONAL COST

## E ANALYSIS OF PROLOG ERRORS

We manually inspected the results in Tab. 1 and identified the following Prolog errors:

- Prolog Parsing Errors:
  - **Prolog query parsing error:** The LLM-generated Prolog query could not be parsed by our Prolog grammar. This error typically occurs due to missing double quotes (see examples in Tab. 10 and 11).
  - **Extraction or verification parsing error:** the LLM-generated answer to the fact extraction or fact verification step is an invalid Prolog literal. This error typically occurs due to nested double quotes (see Tab. 12).
- Prolog Execution Errors:
  - **Intermediate predicate existence error:** Information is missing to solve at least one of the intermediate sub-queries.
  - **Final predicate existence error:** execution reaches the final step, but no match is returned. For this type, we observe two patterns: 1) genuinely missing information in the final sub-query, 2)

Table 7: **Token cost of open-domain QA.** We report mean  $\pm$  1 standard error number of prompt tokens  $\mathbf{P}$  (in thousands) and completion tokens  $\mathbf{C}$  (in thousands) for the results in Tab. 2. We use the vLLM inference engine with prefix caching enabled and report the number of cached tokens in parentheses next to the prompt tokens.

Method	HotpotQA		2WikiMultiHopQA		MuSiQue	
	$P \times 10^3$	$C \times 10^3$	$P \times 10^3$	$C \times 10^3$	$P \times 10^3$	$C \times 10^3$
Standard RAG	3.46 (0.032)	0.159	3.63 (0.032)	0.106	2.21 (0.073)	0.136
Self-Ask	14.5 (10.5)	0.402	15.7 (10.9)	0.682	11.4 (9.06)	0.750
IRCot	62.3 (51.5)	0.047	59.2 (44.6)	0.053	45.3 (38.1)	0.057
$\pi$ -CoT (Ours)	18.0 (5.07)	0.483	21.1 (4.21)	0.543	14.4 (3.28)	0.654

972  
 973 Table 8: **Computational cost of in-context QA.** We report the mean  $\pm 1$  standard error number of  
 974 prompt tokens  $\mathbf{P}$  (in thousands), completion tokens  $\mathbf{C}$  (in thousands), and LLM calls for the results in  
 975 Tab. 4. We use the vLLM inference engine and report the number of cached tokens in parentheses  
 976 next to the prompt tokens when prefix caching was enabled.

977 978 979 980 981 982 983 984 985 986 987 988 989 990 991 992 993 994 995 996 997 998 999 1000 1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015 1016 1017 1018 1019 1020 1021 1022 1023 1024 1025	HotpotQA			2WikiMultiHopQA			MuSiQue		
	Method	$\mathbf{P} \times 10^3$	$\mathbf{C} \times 10^3$	Calls	$\mathbf{P} \times 10^3$	$\mathbf{C} \times 10^3$	Calls	$\mathbf{P} \times 10^3$	$\mathbf{C} \times 10^3$
<i>Llama-3.3-70B-Instruct</i>									
CoT	2.68	0.110	1	2.10	0.0845	1	3.43	0.116	1
$\pi$ -CoT	12.4	0.412	4.37	11.0	0.437	4.47	16.1	0.553	4.90
<i>DeepSeek-R1-Distill-Qwen-32B</i>									
CoT	2.77	0.305	1	2.20	0.298	1	3.57	0.520	1
$\pi$ -CoT	11.1	1.66	4.09	9.59	1.55	4.28	12.3	1.95	4.65

(a) Real-world (Wikipedia-based) multi-hop QA datasets

990 991 992 993 994 995 996 997 998 999	PW-S			PW-M		
	Method	$\mathbf{P} \times 10^3$	$\mathbf{C} \times 10^3$	Calls	$\mathbf{P} \times 10^3$	$\mathbf{C} \times 10^3$
<i>Llama-3.3-70B-Instruct</i>						
CoT	8.12 (8.09)	0.400	1	68.7 (68.6)	0.375	1
$\pi$ -CoT (Ours)	174 (173)	1.59	23.9	2800 (2754)	2.6	40
<i>DeepSeek-R1-Distill-Qwen-32B</i>						
CoT	8.30 (8.26)	1.42	1	70.8 (70.6)	1.13	1
$\pi$ -CoT (Ours)	234 (207)	7.28	21.0	1260 (1230)	7.09	16.8

(b) Synthetic multi-hop QA datasets

1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015 1016 1017 1018 1019 1020 1021 1022 1023 1024 1025	Error Type	HotpotQA	2WikiMultiHopQA	MuSiQue
Parsing: Prolog query parsing error		0.8%	0%	0.8%
Parsing: Execution parsing error		0.8%	0%	0.8%
Execution: Intermediate predicate existence error		16.2%	19.8%	37.4%
Execution: Final predicate existence error		9.2%	0.2%	0.6%
Total errors		27.0%	20.0%	39.6%

1010 Table 9: **Percentage breakdown of Prolog errors for the results in Tab. 1.** We define the errors in  
 1011 App. E and provide examples for each.

1014 mismatches due to Prolog’s strict string equality. The latter can often be resolved when facts  
 1015 are provided back to the model for CoT reasoning. We provide two examples for this type (see  
 1016 examples in Tab. 13 and Tab. 14)

1019 Table 10: **Prolog query parsing error from HotpotQA.** The constant “...Ready for It?” is missing  
 1020 double quotes in the Prolog query.

1022 1023 1024 1025	Question	...Ready for It? is a Taylor Swift song from the album scheduled to be released on what date?
	Prolog Query	album_of_song(...Ready for It?, A1), release_date(A1, A2)

1026  
1027  
1028  
1029Table 11: **Prolog query parsing error from MuSiQue.** Negations (\+) are not allowed by our Prolog parser.1030  
1031  
1032  
1033  
1034  
1035

<b>Question</b>	which professional sports team would you not see play a home game in the arena where the last place Cream performed?
<b>Prolog Query</b>	last_performance_venue("Cream", A1), all_professional_sports_teams(A3), \+home_teams(A1, A3)

1036  
1037  
1038  
1039  
1040Table 12: **Fact extraction parsing error from HotpotQA.** The double quotes in the generated answer are incorrectly escaped.1041  
1042  
1043  
1044  
1045  
1046  
1047

<b>Question</b>	What Cantonese slang term can mean both “ghost man” and to refer to Westerners?
<b>Response</b>	“Gweilo” or “gwallou”

1048  
1049  
1050  
1051Table 13: **Execution error from HotpotQA.** The generated Prolog query involves the sub-query (A1 == A2) and cannot be satisfied by the facts in the knowledge base. Note that the literal “Royal Air Force (RAF)” and the literal “No. 11 Group RAG” are semantically similar, but are not the same under the Prolog operator ==. The ground-truth answer is Royal Air Force.1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062

<b>Question</b>	What where both Hawker Hurricane and No. 1455 Flight apart of?
<b>Prolog Query</b>	part_of("Hawker Hurricane", A1), part_of("No. 1455 Flight", A2), (A1 == A2)
<b>Facts</b>	part_of("Hawker Hurricane", "Royal Air Force (RAF)") part_of("Hawker Hurricane", "Royal Yugoslav Air Force (VVKJ)") part_of("Hawker Hurricane", "Royal Canadian Air Force") part_of("No. 1455 Flight", "No. 11 Group RAF")

1063  
1064  
1065  
1066  
1067  
1068Table 14: **Execution error from MuSiQue.** The generated Prolog query involves a string unification (A1 = A2) and cannot be satisfied by the facts in the knowledge base. The ground-truth answer is John D. Loudermilk.1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

<b>Question</b>	Who wrote turn me on, which was performed by the person who also performed Chasing Pirates?
<b>Prolog Query</b>	performer("Chasing Pirates", A1), performer("turn me on", A2), A1 = A2 -> writer("turn me on", A3)
<b>Facts</b>	performer("Chasing Pirates", "Norah Jones") performer("turn me on", "Sean Smith") writer("turn me on", "Greg Lake") writer("turn me on", "Logan Lynn") writer("turn me on", "Joni Mitchell")

1080 **F ADDITIONAL IN-CONTEXT EXPERIMENTS**  
10811082  
1083 **Table 15: Accuracy of CoT with majority voting.** We report mean exact match (EM) and F1 score  
1084 on a subset of the datasets from Tab. 4

# Samples	HotpotQA		2WikiMultiHopQA		MuSiQue		PW-S	
	EM ↑	F1 ↑	EM ↑	F1 ↑	EM ↑	F1 ↑	EM ↑	F1 ↑
<i>Llama-3.3-70B-Instruct</i>								
1	63.0	79.8	76.8	84.2	49.6	64.3	55.73	75.43
2	63.2	80.3	77.0	84.1	48.4	62.7	54.00	73.95
4	63.4	80.5	76.0	83.4	50.2	63.9	53.93	74.39
8	63.2	80.2	76.0	83.3	51.0	63.8	53.13	73.83
<i>DeepSeek-R1-Distill-Qwen-32B</i>								
1	57.0	74.7	75.8	83.4	45.0	56.4	53.67	74.53
2	57.0	74.8	75.6	83.2	45.6	56.2	54.73	76.14
4	57.0	75.1	76.6	83.4	46.8	57.2	54.47	75.49
8	58.0	75.9	77.4	84.0	49.4	59.3	54.93	75.58

1101 We compare standard CoT to the simplest inference-time intervention method: CoT with majority  
1102 voting. We sample 8 evaluations for each question with the hyperparameters in Tab. 16 and  
1103 compute the performance when majority voting<sup>10</sup> across 2, 4, and 8 of these samples. On  
1104 HotpotQA, 2WikiMultiHopQA, and MuSiQue, CoT with majority voting over 4 samples is similar in  
1105 computational cost to  $\pi$ -CoT (see Tab. 8). Tab. 15 shows that the benefits of repeated sampling over  
1106 single sampling are marginal to none.

1107  
1108 **Table 16: Sampling hyperparameters for the results in Tab. 15.**

Model	Temperature	Top-p	Repetition Penalty	Max Output Tokens
Llama-3.3-70B-Instruct	0.6	0.9	1.0	4096
DeepSeek-R1-Distill-Qwen32B	0.6	0.95	1.0	16384

1114 **G EXAMPLES OF GENERATED PROLOG QUERIES AND DEFINITIONS**  
1115

1116 We show five examples of generated Prolog queries and definitions on HotpotQA, 2WikiMultiHopQA,  
1117 MuSiQue, and PhantomWiki. Examples are randomly chosen.

1120 **H EXAMPLES OF EXECUTION TRACES OF  $\pi$ -CoT**  
1121

1122 We show the full  $\pi$ -CoT workflow on one example from each of 2WikiMultiHopQA, and MuSiQue.  
1123 We use the generated predictions from the experiment of Tab. 1. Examples are randomly chosen.

1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133 <sup>10</sup>We use the majority definition from <https://github.com/EleutherAI/lm-evaluation-harness>.

1134  
1135  
1136  
1137

1138 1139	Question	Query	Target	Definitions
1140 1141 1142 1143 1144 1145	At what age did Cieli di Toscana's singer become blind?	singer("Cieli di Toscana", A1), age_became_blind(A1, A2)	A2	singer(<literal1>, <answer>) -> The singer of <literal1> is <answer>. ; Who is the singer of <literal1>? age_became_blind(<literal1>, <answer>) -> <literal1> became blind at the age of <answer>. ; At what age did <literal1> become blind?
1146 1147 1148 1149 1150 1151 1152 1153	Cooperative Living Organization is located in a city that is the county seat of what county in Florida?	city("Cooperative Living Organization", A1), county_seat(A1, A2)	A2	- city(<literal1>, <answer>) -> The city where <literal1> is located is <answer>. ; What city is <literal1> located in? - county_seat(<literal1>, <answer>) -> <literal1> is the county seat of <answer>. ; What county is <literal1> the county seat of?
1154 1155 1156 1157 1158 1159 1160 1161 1162 1163 1164	Repertoire of Pluecido Domingo appeared in an opera written by whom which was based on a play by William Shakespeare?	opera("Repertoire of Pluecido Domingo", A1), writer(A1, A2), based_on(A1, A3)	A2	- opera(<literal1>, <answer>) -> The opera in which <literal1> appeared is <answer>. ; What opera did <literal1> appear in? - writer(<literal1>, <answer>) -> The writer of the opera <literal1> is <answer>. ; Who wrote the opera <literal1>? - based_on(<literal1>, <answer>) -> The opera <literal1> was based on a play by William Shakespeare <answer>. ; What play by William Shakespeare was the opera <literal1> based on?
1165 1166 1167 1168 1169 1170 1171 1172 1173 1174	In what 2 countries did both X-Large and A Bathing Ape have stores?	x_large_countries(A1), bape_countries(A2), intersection(A1, A2, A3)	A3	- x_large_countries(<answer>) -> X-Large had stores in <answer> countries. ; In what countries did X-Large have stores? - bape_countries(<answer>) -> A Bathing Ape had stores in <answer> countries. ; In what countries did A Bathing Ape have stores? - intersection(<literal1>, <literal2>, <answer>) -> The common countries between <literal1> and <literal2> are <answer>. ; What are the common countries between <literal1> and <literal2>?
1175 1176 1177 1178 1179 1180 1181 1182	Edmund Robert Harris is the principal benefactor of a museum located where ?	museum("Edmund Robert Harris", A1), location(A1, A2)	A2	- museum(<literal1>, <answer>) -> The museum that <literal1> is the principal benefactor of is <answer>. ; What museum is <literal1> the principal benefactor of? - location(<literal1>, <answer>) -> The location of <literal1> is <answer>. ; Where is <literal1> located?

1183  
1184  
1185  
1186  
1187

Table 17: 5 examples on generated queries and definitions on HotpotQA.

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

Question	Query	Target	Definitions
Which film has the director born later, Brutti Di Notte or Bir Türk'E Gönül Verdim?	director("Brutti Di Notte", A1), director("Bir Türk'E Gönül Verdim", A2), date_of_birth(A1, A3), date_of_birth(A2, A4), A3 @> A4 → A5	A5	<ul style="list-style-type: none"> <li>- director(&lt;literal&gt;, &lt;answer&gt;) -&gt; The director of &lt;literal&gt; is &lt;answer&gt;.</li> <li>; Who is the director of &lt;literal&gt;?</li> <li>- date_of_birth(&lt;literal&gt;, &lt;answer&gt;) -&gt; The date of birth of &lt;literal&gt; is &lt;answer&gt;.</li> <li>; When was &lt;literal&gt; born?</li> </ul>
Where did M. K. Muthu's father die?	father("M. K. Muthu", A1), place_of_death(A1, A2)	A2	<ul style="list-style-type: none"> <li>- father(&lt;literal&gt;, &lt;answer&gt;) -&gt; The father of &lt;literal&gt; is &lt;answer&gt;.</li> <li>; Who is the father of &lt;literal&gt;?</li> <li>- place_of_death(&lt;literal&gt;, &lt;answer&gt;) -&gt; The place of death of &lt;literal&gt; is &lt;answer&gt;.</li> <li>; Where did &lt;literal&gt; die?</li> </ul>
What nationality is Julia Parker (Astrologer)'s husband?	husband("Julia Parker (Astrologer)", A1), nationality(A1, A2)	A2	<ul style="list-style-type: none"> <li>- husband(&lt;literal&gt;, &lt;answer&gt;) -&gt; The husband of &lt;literal&gt; is &lt;answer&gt;.</li> <li>; Who is the husband of &lt;literal&gt;?</li> <li>- nationality(&lt;literal&gt;, &lt;answer&gt;) -&gt; The nationality of &lt;literal&gt; is &lt;answer&gt;.</li> <li>; What is the nationality of &lt;literal&gt;?</li> </ul>
Are the directors of both films <i>The Snake Brothers</i> and <i>Kooky</i> from the same country?	director("The Snake Brothers", A1) director("Kooky", A2), country_of_citizenship(A1, A3), country_of_citizenship(A2, A4), A3 == A4 → A5	A5	<ul style="list-style-type: none"> <li>- director(&lt;literal&gt;, &lt;answer&gt;) -&gt; The director of &lt;literal&gt; is &lt;answer&gt;.</li> <li>; Who is the director of &lt;literal&gt;?</li> <li>- country_of_citizenship(&lt;literal&gt;, &lt;answer&gt;) -&gt; The country of citizenship of &lt;literal&gt; is &lt;answer&gt;.</li> <li>; What is the country of citizenship of &lt;literal&gt;?</li> </ul>
What is the place of birth of the director of film <i>Kiss The Sky (Film)</i> ?	director("Kiss The Sky (Film)", A1), place_of_birth(A1, A2)	A2	<ul style="list-style-type: none"> <li>- director(&lt;literal&gt;, &lt;answer&gt;) -&gt; The director of &lt;literal&gt; is &lt;answer&gt;.</li> <li>; Who is the director of &lt;literal&gt;?</li> <li>- place_of_birth(&lt;literal&gt;, &lt;answer&gt;) -&gt; The place of birth of &lt;literal&gt; is &lt;answer&gt;.</li> <li>; Where was &lt;literal&gt; born?</li> </ul>

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

Table 18: 5 examples of generated queries and definitions on 2WikiMultiHopQA.

1242  
1243  
1244  
1245

1246	Question	Query	Target	Definitions
1247	Who was the mother of the person who found the sacred writings that became the Book of Mormon?	founder_of_book_of_mormon(A1), mother(A1, A2)	A2	<ul style="list-style-type: none"> <li>- founder_of_book_of_mormon(&lt;answer&gt;) -&gt; The person who found the sacred writings that became the Book of Mormon is &lt;answer&gt;. ; Who found the sacred writings that became the Book of Mormon?</li> <li>- mother(&lt;literal&gt;, &lt;answer&gt;) -&gt; The mother of &lt;literal&gt; is &lt;answer&gt;. ; Who is the mother of &lt;literal&gt;?</li> </ul>
1248	Who wrote Turn Me On by the Thinking About You performer?	performer("Thinking About You", A1), writer("Turn Me On", A1, A2)	A2	<ul style="list-style-type: none"> <li>- performer(&lt;literal&gt;, &lt;answer&gt;) -&gt; The performer of &lt;literal&gt; is &lt;answer&gt;. ; Who performed &lt;literal&gt;?</li> <li>- writer(&lt;literal1&gt;, &lt;literal2&gt;, &lt;answer&gt;) -&gt; The writer of &lt;literal1&gt; by &lt;literal2&gt; is &lt;answer&gt;. ; Who wrote &lt;literal1&gt; by &lt;literal2&gt;?</li> </ul>
1249	What did M. King Hubbert's employer announce it was in the process of doing in April 2010?	employer("M. King Hubbert", A1), announcement(A1, "April 2010", A2)	A2	<ul style="list-style-type: none"> <li>- employer(&lt;literal&gt;, &lt;answer&gt;) -&gt; The employer of &lt;literal&gt; is &lt;answer&gt;. ; Who is the employer of &lt;literal&gt;?</li> <li>- announcement(&lt;literal&gt;, &lt;date&gt;, &lt;answer&gt;) -&gt; &lt;literal&gt; announced it was in the process of doing &lt;answer&gt; on &lt;date&gt;. ; What did &lt;literal&gt; announce it was doing on &lt;date&gt;?</li> </ul>
1250	What is the experimental satellite being forerunner to communication satellite of INSAT-4CR's manufacturer called?	manufacturer("INSAT-4CR", A1), experimental_satellite(A1, A2)	A2	<ul style="list-style-type: none"> <li>- manufacturer(&lt;literal&gt;, &lt;answer&gt;) -&gt; The manufacturer of &lt;literal&gt; is &lt;answer&gt;. ; Who is the manufacturer of &lt;literal&gt;?</li> <li>- experimental_satellite(&lt;literal&gt;, &lt;answer&gt;) -&gt; The experimental satellite being the forerunner to the communication satellite of &lt;literal&gt; is &lt;answer&gt;. ; What is the experimental satellite being the forerunner to the communication satellite of &lt;literal&gt;?</li> </ul>
1251	The Socialist Autonomous Province of the city where there were mass executions of Danube Swabian population are located in where?	city_with_mass_executions("Danube Swabian", A1), socialist_autonomous_province(A1, A2), location(A2, A3)	A3	<ul style="list-style-type: none"> <li>- city_with_mass_executions(&lt;literal&gt;, &lt;answer&gt;) -&gt; The city with mass executions of &lt;literal&gt; population is &lt;answer&gt;. ; What city had mass executions of the &lt;literal&gt; population?</li> <li>- socialist_autonomous_province(&lt;literal&gt;, &lt;answer&gt;) -&gt; The Socialist Autonomous Province of &lt;literal&gt; is &lt;answer&gt;. ; What is the Socialist Autonomous Province of &lt;literal&gt;?</li> <li>- location(&lt;literal&gt;, &lt;answer&gt;) -&gt; &lt;literal&gt; is located in &lt;answer&gt;. ; Where is &lt;literal&gt; located?</li> </ul>

1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

Table 19: 5 examples on generated queries and definitions on MuSiQue.

1296

1297

1298

1299

1300

1301

Question	Query	Target	Definitions
How many children does Shelly Reece have?	aggregate_all(count, distinct(child("Shelly Reece", A1)), A2)	A2	<ul style="list-style-type: none"> <li>- child(&lt;literal&gt;, &lt;answer&gt;) -&gt; The child of &lt;literal&gt; is &lt;answer&gt;. ; Who is the child of &lt;literal&gt;?</li> </ul>
How many friends does the friend of the person whose hobby is aerospace have?	hobby(A1, "aerospace"), friend(A1, A2), aggregate_all(count, distinct(friend(A2, A3)), A4)	A4	<ul style="list-style-type: none"> <li>- hobby(&lt;answer&gt;, &lt;literal&gt;) -&gt; The hobby of &lt;answer&gt; is &lt;literal&gt;. ; Who is the person whose hobby is &lt;literal&gt;?</li> <li>- friend(&lt;literal&gt;, &lt;answer&gt;) -&gt; The friend of &lt;literal&gt; is &lt;answer&gt;. ; Who is the friend of &lt;literal&gt;?</li> </ul>
How many children does the mother of the child of the person whose date of birth is 0985-04-02 have?	dob(A1, "0985-04-02"), child(A1, A2), mother(A2, A3), aggregate_all(count, distinct(child(A3, A4)), A5)	A5	<ul style="list-style-type: none"> <li>- dob(&lt;answer&gt;, &lt;literal&gt;) -&gt; The date of birth of &lt;answer&gt; is &lt;literal&gt;. ; Who is the person whose date of birth is &lt;literal&gt;?</li> <li>- child(&lt;literal&gt;, &lt;answer&gt;) -&gt; The child of &lt;literal&gt; is &lt;answer&gt;. ; Who is the child of &lt;literal&gt;?</li> <li>- mother(&lt;literal&gt;, &lt;answer&gt;) -&gt; The mother of &lt;literal&gt; is &lt;answer&gt;. ; Who is the mother of &lt;literal&gt;?</li> </ul>
How many mothers does the friend of the sister of the sister of the child of the parent of the person whose hobby is mineral collecting have?	hobby(A1, "mineral collecting"), parent(A1, A2), child(A2, A3), sister(A3, A4), sister(A4, A5), friend(A5, A6), aggregate_all(count, distinct(mother(A6, A7)), A8)	A8	<ul style="list-style-type: none"> <li>- hobby(&lt;answer&gt;, &lt;literal&gt;) -&gt; The hobby of &lt;answer&gt; is &lt;literal&gt;. ; Who is the person whose hobby is &lt;literal&gt;?</li> <li>- parent(&lt;literal&gt;, &lt;answer&gt;) -&gt; The parent of &lt;literal&gt; is &lt;answer&gt;. ; Who is the parent of &lt;literal&gt;?</li> <li>- child(&lt;literal&gt;, &lt;answer&gt;) -&gt; The child of &lt;literal&gt; is &lt;answer&gt;. ; Who is the child of &lt;literal&gt;?</li> <li>- sister(&lt;literal&gt;, &lt;answer&gt;) -&gt; The sister of &lt;literal&gt; is &lt;answer&gt;. ; Who is the sister of &lt;literal&gt;?</li> <li>- friend(&lt;literal&gt;, &lt;answer&gt;) -&gt; The friend of &lt;literal&gt; is &lt;answer&gt;. ; Who is the friend of &lt;literal&gt;?</li> <li>- mother(&lt;literal&gt;, &lt;answer&gt;) -&gt; The mother of &lt;literal&gt; is &lt;answer&gt;. ; Who is the mother of &lt;literal&gt;?</li> </ul>
How many children does the sibling of the brother of the brother of the sister of Jamal Song have?	sister("Jamal Song", A1), brother(A1, A2), brother(A2, A3), sibling(A3, A4), aggregate_all(count, distinct(child(A4, A5)), A6)	A6	<ul style="list-style-type: none"> <li>- sister(&lt;literal&gt;, &lt;answer&gt;) -&gt; The sister of &lt;literal&gt; is &lt;answer&gt;. ; Who is the sister of &lt;literal&gt;?</li> <li>- brother(&lt;literal&gt;, &lt;answer&gt;) -&gt; The brother of &lt;literal&gt; is &lt;answer&gt;. ; Who is the brother of &lt;literal&gt;?</li> <li>- sibling(&lt;literal&gt;, &lt;answer&gt;) -&gt; The sibling of &lt;literal&gt; is &lt;answer&gt;. ; Who is the sibling of &lt;literal&gt;?</li> <li>- child(&lt;literal&gt;, &lt;answer&gt;) -&gt; The child of &lt;literal&gt; is &lt;answer&gt;. ; Who is the child of &lt;literal&gt;?</li> </ul>

1345

1346

1347

1348

1349

Table 20: 5 examples on generated queries and definitions on PhantomWiki-S.

---

1350  
1351     **Question**  
1352     Did John Updike and Tom Clancy both publish more than 15 bestselling novels?  
1353     **Query**  
1354     bestselling\_novels(John Updike, A1),  
1355     bestselling\_novels("Tom Clancy", A2),  
1356     (A1 @> 15 -> B1 = "Yes" ; B1 = "No"),  
1357     (A2 @> 15 -> B2 = "Yes" ; B2 = "No"),  
1358     (B1 == "Yes" -> (B2 == "Yes" -> A3 = "Yes" ; A3 = "No") ; A3 = "No")

---

1359  
1360     **Definitions**  
1361     bestselling\_novels(<literal1>, <answer>) → <literal1> published <answer> bestselling  
1362     novels. How many bestselling novels did <literal1> publish?

---

1363     **Sub-Query 1:** How many bestselling novels did John Updike publish?  
1364     **Retrieved passages:** *John Updike bibliography, Tom Wolfe, Toward the End of Time, John Updike, John Updike*  
1365     **Answer:** more than 20  
1366     **Sub-Query 2:** How many bestselling novels did Tom Clancy publish?  
1367     **Retrieved passages:** *Duty and Honor, Line of Sight (novel), Tom Clancy's EndWar, Mike Maden, Mark Greaney (novelist)*  
1368     **Answer:** 17

---

1369  
1370     **Notes**  
1371     John Updike published more than 20 bestselling novels.  
1372     Tom Clancy published 17 bestselling novels.

1373  
1374     **Prolog Answer:**  
1375     **Final Answer:** Yes

---

1376  
1377     Table 21: One example on the execution trace of a question on HotpotQA.

---

1378  
1379     **Question**  
1380     Who is Wisigard's father-in-law?  
1381     **Query**  
1382     spouse(Wisigard, A1),  
1383     father(A1, A2)

---

1384     **Definitions**  
1385     spouse(<literal>, <answer>)^ -> The spouse of <literal> is <answer>. ; Who is the  
1386     spouse of <literal>? father(<literal>, <answer>)^ -> The father of <literal> is <answer>. ; Who is the father of <literal>?

---

1387  
1388     **Sub-Query 1:** Who is the spouse of Wisigard?  
1389     **Retrieved passages:** *Wisigard, Wisigard, Deuteria, Deuteria, Theudebert I*  
1390     **Answer:** Theudebert I

1391  
1392     **Sub-Query 2:** Who is the father of Theudebert I?  
1393     **Retrieved passages:** *Theudebert of Soissons, Theudebert of Soissons, Theudebert I, Theudebert I, Theudebert I*  
1394     **Answer:** Theuderic I

---

1395  
1396     **Notes**  
1397     The spouse of Wisigard is Theudebert I.  
1398     The father of Theudebert I is Theuderic I.

1399  
1400     **Prolog Answer:** Theuderic I  
1401     **Final Answer:** Theuderic I

---

1402  
1403     Table 22: One example on the execution trace of a question on 2WikiMultiHopQA.

1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413

---

**Question**

In what region of the country containing A Luoi is the birthplace of John Phan located?

---

**Query**

contains("A Luoi", A1),  
birthplace("John Phan", A2),  
region(A1, A2, A3)

---

**Definitions**

contains(<literal>, <answer>)<sup>‘</sup> -> <literal> is contained in the country <answer>. ;  
What country contains <literal>?  
birthplace(<literal>, <answer>)<sup>‘</sup> -> The birthplace of <literal> is <answer>. ; What is  
the birthplace of <literal>?  
region(<literal1>, <literal2>, <answer>)<sup>‘</sup> -> <literal2> is located in the <answer>  
region of <literal1>. ; What region of <literal1> is <literal2> located in?

---

**Sub-Query 1:** What country contains A Luoi?

**Retrieved passages:** *A Luoi Camp, A Luoi District, A Luoi District, A Luoi Camp, A Luoi District*

**Answer:** Vietnam

**Sub-Query 2:** What is the birthplace of "John Phan"?

**Retrieved passages:** *John Phan, Phan Boi Châu, Peter C. Phan, John Phan, Phan Đình Phùng*

**Answer:** Da Nang, Vietnam

**Sub-Query 3:** What region of "Vietnam" is "Da Nang, Vietnam" located in?

**Retrieved passages:** *Da Nang, Da Nang Air Base, Da Nang University of Economics, Hoàng Sa District, Da Nang*

**Answer:** South Central Coast

---

**Notes**

A Luoi is contained in the country Vietnam.  
The birthplace of John Phan is Da Nang, Vietnam.  
Da Nang, Vietnam is located in the South Central Coast region of Vietnam.

**Prolog Answer:** South Central Coast

---

**Final Answer:** South Central Coast

---

1448 Table 23: 1 example on the execution trace of a question on MuSiQue.  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457