Large Language Models are In-Context Semantic Reasoners rather than Symbolic Reasoners

Anonymous Author(s)

Affiliation Address email

Abstract

The emergent few-shot reasoning capabilities of Large Language Models (LLMs) have excited the natural language and machine learning community over recent years. Despite of numerous successful applications, the underlying mechanism of such in-context capabilities still remains unclear. In this work, we hypothesize that the learned semantics of language tokens do the most heavy lifting during the reasoning process. Different from human's symbolic reasoning process, the semantic representations of LLMs could create strong connections among tokens, thus composing a superficial logical chain. To test our hypothesis, we decouple semantics from the language reasoning process and evaluate three kinds of reasoning abilities, i.e., deduction, induction and abduction. Our findings reveal that semantics play a vital role in LLMs' in-context reasoning—LLMs perform significantly better when semantics are consistent with commonsense but struggle to solve symbolic or counter-commonsense reasoning tasks by leveraging in-context new knowledge. The surprising observations question whether modern LLMs have mastered the inductive, deductive and abductive reasoning abilities as in human intelligence, and motivate research on unveiling the magic existing within the black-box LLMs. On the whole, our analysis provides a novel perspective on the role of semantics in developing and evaluating language models' reasoning abilities.

1 Introduction

2

3

8

9

10

12

13

14

15

16

17

18

19

In recent years, Large Language Models (LLMs) have achieved impressive performance on a variety of natural language tasks, including question answering, text summarization, machine translation, logic reasoning, *etc*. These successes have been largely attributed to the emergent ability of LLMs to utilize a "zero-shot" or "few-shot" learning approach without any gradient updates—a task description or a few examples are provided to guide their reasoning process [1–4]. One typical example is the "chain-of-thought (CoT)" approach, involving reasoning demonstrations or a simple prompt such as "Let's think step by step" to perform complex reasoning tasks [5, 6].

Despite the powerful and versatile in-context learning ability of LLMs, the underlying mechanisms by which they operate within a given context still remain unclear. Previous works investigate which aspects of the given examples contribute to the final task performance, including ground-truth labels and example ordering [7–9]. Another line of recent work has focused on explaining and leveraging the in-context learning (ICL) mechanism [10–13]. However, the basic problem they have in common is that the in-context prompts they input are based on natural language queries to investigate the reasoning abilities of LLMs. According to the Dual Process Theory [14, 15], humans are capable of using symbolic reasoning with System II to solve complex logical reasoning problems. To fill the research gap, we systematically study the in-context reasoning ability of LLMs by decoupling the

semantics from the language reasoning process. With extensive experiments, we aim to answer the following research question: *Are LLMs good in-context reasoners without semantics?*

In this work, we hypothesize that the learned semantics of language tokens play an important role in 38 the reasoning process, creating strong connections among tokens which help to compose a superficial 39 logical chain (shortcut) instead of really performing the formal reasoning process. To test our 40 hypothesis, given symbolic knowledge (facts and rules), we test three kinds of reasoning abilities (i.e., deduction, induction, abduction) on a newly proposed synthetic dataset: Symbolic Tree dataset, which is composed of closed-world, noise-free, multi-hop symbolic reasoning data generated with 43 logical rules. Besides, we also experiment with ProofWriter [16] task, containing questions whose 44 answers require multi-hop reasoning. Our findings suggest that semantics indeed play a vital role 45 in LLMs' in-context reasoning: When semantics are consistent with commonsense, LLMs perform 46 fairly well; when semantics are decoupled or in the counter-commonsense context, LLMs struggle to solve the reasoning tasks by leveraging in-context new knowledge. Moreover, we also study the 48 memorization ability of LLMs to memorize new symbols and semantics information, allowing us to investigate the role of semantics on LLMs' knowledge update ability.

To the best of our knowledge, this is the first study of the effect of semantics on LLMs' in-context reasoning abilities. Our analysis underscores the importance of semantics in LLMs' reasoning ability and questions whether modern LLMs have mastered the formal reasoning abilities as in human intelligence. We hope our findings can provide a novel perspective on the role of semantics in LLMs' in-context abilities, and inspire further research on unveiling the magic inside the black-box LLMs.

56 2 Related Works

57

58

59

61

62

63

64

65

66

67

68

70

71

72

73

74 75

76 77

78

79

80

81

82

83

85

86

87

88

Reasoning in LLMs Reasoning is a fundamental cognitive process involving logical inferences and conclusions based on given information. Developing models with strong reasoning capabilities has attracted increasing attention [17]. Since then, various benchmarks focusing on different aspects of reasoning have been proposed, including natural language inference (NLI) [18-20], commonsense reasoning [21, 22], multi-hop reasoning [23, 24], logic reasoning [25–27] etc. These benchmarks are evaluated within semantics. However, we focus on decoupling semantics to test LLMs' pure symbolic logic reasoning capabilities. In recent years, there has been growing interests in studying the reasoning abilities of LLMs. For example, "chain-of-thought (CoT)" [5, 28, 6] is proposed to facilitate models to generate a reasoning path that decomposes complex reasoning into multiple easier steps. Creswell et al. [29] solve multi-step reasoning tasks by interacting between selection and inference to generate immediate reasoning steps. This significantly improves the performance on arithmetic [30], commonsense [21, 31], and symbolic reasoning [5] benchmarks. There has been also some work to explore the effect of semantics on reasoning. For example, Dasgupta et al. [32] evaluate three logical reasoning tasks, namely natural language inference (NLI), syllogisms and Wason selection based on whether the content of the problem is aligned with prior knowledge, concluding that LLMs show human-like content effects on reasoning. Schlegel et al. [33] also reach the similar conclusion. However, this is just a subordinate conclusion (LLMs are in-context semantic reasoners that rely on shallow semantic associations for prediction; consequently their predictions are naturally influenced by semantics). It cannot be used to extrapolate the conclusion that "LLMs are not symbolic reasoners. In contrast, our task focuses on pure symbolic reasoning. When we feed these to LLMs, they show significantly worse performance compared to when normal semantic words are fed. This phenomenon indicates that LLMs fail to invoke the basic formal reasoning abilities of humans but instead rely on shallow semantic associations for prediction.

In-Context Learning LLMs' reasoning abilities are closely related to in-context learning (ICL). ICL refers to the ability of language models to adapt and learn from a few prompt examples during the inference process. In recent years, there has been a focus on exploring how to improve the performance of ICL [34–38]. Despite the empirical success, the underlying mechanisms of ICL still remain unclear. A few studies have shown how the chain-of-thought prompting effect the reasoning abilities of LLMs and which aspects of the prompting contribute to their performance [8, 39]. Specifically, the order of demonstrations may lead to large performance fluctuations [40, 9]. Wang et al. [41] explores factors such as the order of reasoning steps, whether the demonstrations consist of invalid reasoning steps or are related to the query. Turpin et al. [42] introduces bias features into CoT demonstrations, highlighting the potential for misleading explanations even within a chain of thought. Madaan and Yazdanbakhsh [43] aims to gain a deeper understanding of which aspects of CoT-based

prompting affect LLMs' reasoning performance by conducting counterfactual experiments. Recent works also explore the effect of ground-truth labels and question the necessity of ground-truth input-output mapping—using incorrect labels in the examples only marginally lowers the performance [44] and input-label correspondence plays a more important role in contextual demonstration [45]. Wei et al. [46] investigate the effects of semantic priors and input-label mapping on in-context learning using different-scale and instruction-tuned models. Wei et al. [47] propose symbol tuning, which involves fine-tuning a language model using input-label mapping unrelated to semantics prior in order to investigate whether LLMs can induce input-label patterns further improve reasoning abilities. To further understand why in-context learning works, some work provides theoretical analysis that in-context learning can be formalized as Bayesian inference [13] or some instances of ICL can be understood as implicit implementation of known learning algorithms [48]. However, the existing analyses of ICL are mainly based on natural language input with rich semantic information. We hypothesize that this might not be able to reflect their true level of symbolic reasoning abilities, including deduction, induction and abduction. Thus, this paper aims to decouple semantics in LLMs' in-context reasoning abilities.

Symbolic Reasoning Symbolic reasoning has long been studied in the field of artificial intelligence [49-53] and cognitive science [54-58]. It involves manipulating symbols and applying logical rules to perform deduction [59], induction [55], and abduction [60]. Recently, there has been some work to explore LLMs' ability of symbolic reasoning. Qian et al. [61] evaluate a set of simple symbolic manipulation (copy, reverse and addition) tasks to uncover the difficulty of the LLMs in handling OOD symbolic generalization, highlighting the limitation in arithmetic and symbolic induction. BIG-Bench [62] contain the symbol_interpretation task aimed at reasoning and interpreting a simple scene consisting of some defined objects. Shin et al. [63] explore using LLM-based models for program synthesis by leveraging inferred execution traces to guide the generation of correct programs. Lample and Charton [64] focus on applying LLM-based models to mathematical reasoning, proposing a framework that combines deep learning with symbolic mathematics to perform algebraic reasoning, equation solving, and theorem proving. Pallagani et al. [65] use LLMs for automated planning—a branch of AI concerned with realizing action sequences (plans) to achieve certain goals, typically executed by intelligent agents, autonomous robots, and unmanned vehicles. Gao et al. [66] use LLMs to read natural language problems and generate programs as the intermediate reasoning steps proposes a new method using a symbolic module to address to improve LLMs' symbolic reasoning abilities.

3 Decoupling Semantics from In-Context Reasoning

124 3.1 Task Definitions

To begin, we first introduce the definitions of reasoning and memorization mechanisms and provide task descriptions for each. Examples of the tasks are shown in Figure 1.

Reasoning In the field of psychology, reasoning refers to the process of using logical operations to draw conclusions or make inferences based on available information [67–70]. As an abstract notion, it encompasses a variety of aspects. Traditionally, we can classify it into three categories:

- Deductive reasoning is a logical process in which a conclusion can be derived from given premises or principles, meaning predicting new facts based on existing facts and logical rules. For example, given the two facts (Lisa, sisterOf, Alice) and (Alice, motherOf, Bob) along with a logical rule ∀x, y, z : sisterOf(x, y) ∧ motherOf(y, z) → auntOf(x, z), the new fact (Lisa, auntOf, Bob) can be derived through deductive reasoning. The task is to predict the True/False of a predicted fact given facts and rules. The accuracy is the proportion of correct predictions.
- Inductive reasoning involves making generalizations based on specific observations or evidence. In other words, a logical rule can be induced from given facts. For instance, given a set of observations that person A is the parent of person B and person B is the child of person A, inductive reasoning is to conclude the logical rule $\forall x, y : \operatorname{parentOf}(x, y) \to \operatorname{childOf}(y, x)$. We perform the rule generation task. Given multiple facts with similar patterns and a rule template, the goal is to induce a rule that entails these facts. We test the generated rules against the ground truth rules. If the generated rule matches the ground truth rule exactly, we regard the rule as correct; otherwise, we regard the rule as incorrect. The precision is the proportion of correct rules. More details of the rule template and the ground-truth rules are provided in Appendix G.

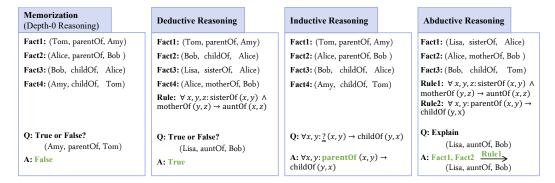


Figure 1: Task Definitions. **Memorization**: retrieving the predicted fact from in-context knowledge. **Deductive:** predicting the correctness of the predicted fact given rules and facts. **Inductive:** generating a rule based on multiple facts with similar patterns. Abductive: explaining the predicted fact based on given rules and facts.

• Abductive reasoning is a logical process of seeking a hypothesis that best fits or explains a set of observations. For example, given a lot of facts including (Lisa, sisterOf, Alice) and (Alice, motherOf, Bob), along with a set of logical rules including $\forall x, y, z$: sisterOf $(x, y) \land$ motherOf $(y, z) \rightarrow \text{auntOf}(x, z)$, if we observe Lisa is Bob's aunt, one possible explanation is that Lisa is Alice's sister and Alice is Bob's mother. We use explanation generation to evaluate the abductive reasoning ability. Given a theory including facts and logical rules, the task is to select specific facts and a logical rule from the given theory to explain the observation. The observation is chosen from inferred facts. We use Proof Accuracy (PA) as an evaluation metric, *i.e.*, the fraction of examples where the generated proof matches exactly any of the gold proofs.

Memorization Memory plays a crucial role in reasoning, as it involves storing the in-context or parametric knowledge necessary for the reasoning process. In some sense, memory can be considered as Depth=0 reasoning, where the question is a known fact. The reasoning task involves retrieving the fact itself from in-context or knowledge within the parameters. However, the specific impact of semantics on memorization has not been extensively explored. Thus, in addition to decoupling semantics from reasoning, we also try to study the impact of semantics on memorization. Specifically, we use a new dataset to fine-tune a language model and test its time, efficiency and forgetting ratio: time is the fine-tuning time cost of adding/updating facts, efficiency is the filter MRR (the mean reciprocal of rank of the correct entity [71]) of the facts added/updated, and forgetting ratio is the filter MRR of the facts that should not be updated. When evaluating whether a fact has been successfully added or updated, we query LLM with a question about the tail entity and rank the probability of the true tail against all entities. The better LLM remembers a triplet, the higher the MRR gets.

Evaluation Datasets

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159 160

161

162

163

164

165

166

171

172

173

174

175

176 177

178

179

181

Our goal is to decouple semantics from the in-context reasoning process and solely rely on the given 167 (new) knowledge to perform reasoning tasks. To implement this, we use Symbolic Tree [72] and 169 ProofWriter [16] datasets, which contain both relevant and irrelevant facts and LLMs need to infer the unknown facts after selecting relevant facts from memory. 170

The Symbolic Tree dataset is an artificially close-world and noise-free symbolic dataset generated with complex logical rules. The dataset consists of randomly sampled "basic facts", which include gender information and "parentOf" relations among individuals. With the given logical rules, the dataset allows for reasoning about 28 different types of family relations, ranging from easy inferences (e.g., fatherhood), to more elaborate ones (e.g., a daughter of someone's cousin). Facts consist of basic facts (in-context knowledge) and inferred facts (what to reason). Note that Symbolic Tree is a close-world dataset, which means that any facts not presented in the dataset are assumed to be false. Thus, we construct the false facts by replacing the head entity or tail entity with a random entity as negative examples in *inferred facts*. Considering the context window size limitation, we restrict each tree's depth to 5 to generate the dataset. We experiment with 10 sampled Symbolic Trees; each has 180 30 kinds of relations (28 inferred relations, gender and parentOf relation), 26 entities, about 35 basic facts, 300 inferred facts and 300 false ones.

To decouple the semantics within the dataset, we replace the relation names (such as "parent") 183 with hand-crafted symbols (e.g., "r1", "r2", ...), so that LLMs cannot leverage the semantics of the 184 predicates in reasoning but must resort to the given new knowledge (presented as in-context facts 185 and rules). We also experiment with replacing entity names (such as "Alice") with "e1", "e2", ..., 186 but find that it has little impact on performance (more details are provided in Appendix Q). During 187 the symbol generation process, we also try to randomly sample some letters as relation names (e.g.,188 189 "lnggv" instead of "r1"), but we observe that LLMs struggle to understand garbled characters, which may negatively affect performance (further discussion is provided in Appendix N). 190

ProofWriter [16] tasks provide artificial facts and rules expressed in natural language. For our 191 experiments, we use a subset of the ProofWriter Open World Assumption (OWA) dataset with a 192 depth of 1, 2, 3 and 5 (there is no depth 4 task), which contains many small rulebases of facts and 193 rules, expressed in English and do not exist in LLMs' knowledge base. Each rulebase has a set of questions (English statements) that can be proven true, false or "Unknown". Note that if we want to 195 prove something Unknown, it is necessary to enumerate all possible facts and check their true/false. Thus, we remove all the Unknowns and replace the subject and object with entity IDs. This dataset is 197 simpler than Symbolic Tree. Considering most of the predicates in the sentences are unmeaningful 198 verbs like "is" and "can", we only replace the entities with entity IDs to decouple semantics. Take 199 "Anne is kind." as an example. We substitute subject (Anne) and object (kind) with "e1" and "e2", 200 respectively, resulting in "e1 is e2". Figure 2 provides an illustrative example. 201

4 Experiment

202

203

204

205

206

207

209

210

211

212

213

214

215

216

217

218

219

220

221

224

225

226

227

228

229

230

231

232

233

234

235

Models Selected for Evaluation We primarily evaluate the performance of ChatGPT, GPT-4 and LLaMA. ChatGPT and GPT-4 are advanced AI models developed by OpenAI and have demonstrated strong reasoning abilities across various tasks and benchmarks. LLaMA is an open-source large language model developed by Meta AI, with number of parameters ranging from 7B to 65B. Due to computational resource constraints, we could only fine-tune the LLaMA-7B version, which is used in our memorization test. Note that in our study, we did attempt reasoning experiments using finetuned LLaMA-7b model. It performs significantly worse in reasoning tasks and even struggles to understand the instructions. Additionally, tasks requiring extensive facts and logical rules are hindered by the limited context window size. As a result, we did not conduct reasoning experiments with it. Additionally, when comparing the Given a set of rules and facts, you have to Given a set of rules and facts, you have reason whether a statement is true or false. to reason whether a statement is true or Here are some facts and rules Here are some facts and rules: The bear likes the dog. The cow is ro The e4 likes the e5. The cow likes the bear The e14 is The e14 likes the e4 The cow needs the bea The e14 needs the e The dog needs the squirrel The squirrel needs the dos The e5 sees the e14. If someone is round then they like the The e26 needs the auirrel. If someone is e2 then they like the e26 If the bear is round and the bear likes the If the e4 is e2 and the e4 likes the e26 squirrel then the squirrel needs the be then the e26 needs the If the cow needs the dog then the cow is If the e14 needs the e5 then the e14 is Does it imply that the statement "The cow Does it imply that the statement "The likes the squirrel." is True? e14 likes the e26." is True?

Figure 2: Decoupling semantics from the ProofWriter task. In the original ProofWriter task, entities are represented by their names (left). However, in our decoupled setting, we replace the entity names with unique entity IDs (right).

reasoning abilities of LLMs, we also use some **logic-based** symbolic methods to conduct experiments as the baseline. To compare memorization, we use a popular graph database **Neo4j** [73] as the baseline. To ensure a relatively fair comparison, we configure Neo4j with a pre-stored knowledge base that has comparable disk space size to LLaMA. More introduction of Neo4j is represented in Appendix F.

Evaluation Setup For reasoning, we use Symbolic Tree and ProofWriter as evaluation data. We refer to the raw data, where semantics are retained, as *Semantics*. When semantics are decoupled using symbols, we refer to it as *Symbols*. For the Symbolic Tree dataset, we experiment with 10 sampled trees and report the average results, where facts and rules can be represented as logical language and natural language text as the input of LLMs. For example, the fact "motherOf(Alice, Bob)" can be represented as "Alice is Bob's mother"; the rule " $\forall x, y$: parentOf(x, y) \rightarrow childOf(x, y)" can be represented as "If x is parent of y, then y is child of x.". Through numerous trials, we find that for the *Symbols* setting, LLMs tend to perform better when using logic language representations. Conversely, for the *Semantics* setting, LLMs tend to perform better when using natural language text. We select the representation that yields better performance in LLMs' reasoning. Additional results are presented

in Appendix M. We consider zero-shot, zero-shot CoT, few-shot CoT and zero-plus-few-shot-CoT as baselines. To generate explanations for few-shot CoT experiments, for deductive reasoning, we use zero-shot CoT (*i.e.*, Let's think step by step) to generate explanations given the random questions; for abductive reasoning, we randomly select five examples and manually design their demonstrations. We provide all prompts and CoT demonstrations in Appendix A. We use the accuracy of various tasks as the reasoning result, including deducing the correctness of a conclusion, inducing correct rules, or finding explanations for hypotheses.

For memorization, we randomly selected 1,258 triplets from four sampled Symbolic Trees to fine-tune 245 the LLaMA. After adding these triplets, we perform a second fine-tuning step where we update 246 half of the added triplets. To obtain the updated facts, we select the triplets in the first two trees 247 and replace the tail entities with other random entities. Since these updates are chosen from two 248 independent Symbolic Trees, they did not overlap with the remaining half of the facts. We then used 249 the other two trees to evaluate the impact of updating knowledge on other knowledge, namely the forgetting ratio. We still use *Symbols* and *Semantics* to denote different experiment settings. Both settings ensure that the new information provided does not overlap with the old knowledge base of LLMs, avoiding any ambiguation problems and eliminating the influence of pre-existing knowledge 253 on the memorization task. When testing, we follow the prompting of Taori et al. [74], using the head 254 entity and relation as instructions and providing all candidate tails as input. The detailed prompts are 255 contained in Appendix A. 256

Implementation Details For ChatGPT and GPT-4, we use the chat completion API provided by OpenAI. We use a temperature of zero to generate output. Additionally, we set the frequency penalty to zero and top p to 1, which are the default values for these APIs.

For LLaMA-7B, we utilized 4 A100 80G GPUs with batch size 64 for finetuning. The training process involved 100 epochs, employing a cosine learning rate schedule with an initial learning rate of 2e-5. We run these experiments three times and recorded their mean MRR and standard deviations. Please refer to Appendix H for more details.

For logic-based symbolic baseline, in the deductive reasoning setting, it enumerates paths between head h and tail t and uses activated rules to infer the answer; For inductive reasoning, we adopt AMIE+ [75], which first enumerates possible rules and then learns a scalar weight for each rule to encode its quality. For abductive reasoning, we locate the logical rule that reason about the relation of the fact and find all paths connecting the head and tail that can activate the rule. These path facts, along with the logical rule, serve as the explanations.

4.1 Semantics Matter in LLMs' memorizing

271

272

283

284

285

286

288

We first test the memorization ability of LLMs when new knowledge are presented in semantics/symbols forms. The results are reported in Table 1.

Results From Table 1, the *efficiency* of adding and updating semantic knowledge is higher compared 273 to symbolic knowledge. This suggests that semantic knowledge is easier for LLMs to memorize than symbolic knowledge, similar to human's memory capabilities (memorizing symbols is generally more challenging than memorizing words with semantic meanings). However, we also find that the forgetting ratio of Semantics setting is higher than the symbolic setting. This could be attributed to the fact that semantic knowledge has stronger correlation with each other than symbolic knowledge 278 in LLMs. In other words, LLMs may utilize shallow semantic associations for memorization. When 279 a portion of knowledge is updated, it can inadvertently affect other knowledge that should remain 280 281 unaffected. In contrast, symbolic LLMs rely on rote memorization, which makes them less susceptible 282 to such inadvertent changes and forgetting.

We also compare fine-tuned language models with the deterministic graph DB Neo4J to explore the memorization abilities of neural-based and symbolic-based methods (More illustrations are included in Appendix E). From the results shown in Table 4.1, we can see that knowledge update using Neo4j achieves 100% accuracy when inserting new triplets or editing existing triplets, regardless of whether the knowledge is symbolic or semantic. As expected, since the added or updated knowledge does not overlap with the existing knowledge base, there is no further influence on the existing knowledge in the database. Additionally, compared to the computational cost of fine-tuning LLMs, updating knowledge in a graph database with optimized storage mechanisms is significantly faster. This affirms

Table 1: Memorization abilities: LLaMA-7B and Neo4j. MRR are in %.

Method	Category	Adding Efficiency (MRR)	Updating Efficiency (MRR)	Forgetting (MRR↓)	Time/1k triplets
LLaMA-7B	Semantics Symbols	$ 50.375 \pm 1.27 48.91 \pm 4.3 $	$51.34 \pm 0.55 40.74 \pm 2.26$	$7.02 \pm 1.55 \\ 2.2 \pm 0.99$	41.5 min 41.5 min
Neo4j	Semantics Symbols	100 100	100 100	0	19s 19s

Table 2: The reasoning results of Symbolic Tree. Results are in %.

Category	Model	Baseline	deduction	induction	abduction
		Zero-Shot	52.6	6.10	1.50
	ChatGPT	Zero-Shot-CoT	55.7	7.86	4.90
	ChaiGFI	Few-Shot-CoT	54.8	-	18.2
Symbols		Zero-Plus-Few-Shot-CoT	55.7	-	16.8
-		Zero-Shot	68.8	9.28	25.0
	GPT-4	Zero-Shot-CoT	71.1	8.93	31.2
		Few-Shot-CoT	67.6	-	44.2
	ChatGPT	Zero-Shot	66.1	36.4	2.94
		Zero-Shot-CoT	65.5	32.2	3.40
		Few-Shot-CoT	67.1	-	21.8
Semantics		Zero-Plus-Few-Shot-CoT	67.2	-	20.9
		Zero-Shot	79.2	52.5	27.3
	GPT-4	Zero-Shot-CoT	86.2	53.9	33.4
		Few-Shot-CoT	91.1	-	69.2
	Random	-	50.1	3.57	-
	Logic-based	-	100	57.1	100

the huge advantage of using KGs/external DBs to update knowledge rather than finetuning, aligning with the recent trend of retrieval-based LLMs.

4.2 Semantics Play a Vital Role in LLMs' Reasoning

293

297

298

299

300

301

In this section, we evaluate the impact of decoupling semantics from LLMs' in-context reasoning. In Table 2, we present the results of deductive, inductive, and abductive reasoning tasks on the Symbolic Tree datasets.

Results From Table 2, we observe that in all reasoning scenarios, *Semantics* setting significantly outperforms *Symbols* setting. Notably, in the inductive experiments, *Semantics* achieves approximately 30% higher absolute accuracy compared to *Symbols* setting. This indicates that preserving rich semantics in the reasoning process leads to better performance for LLMs.

Despite the improved in-context reasoning performance of LLMs with rich semantics, when compared to logic-based symbolic methods, LLMs still exhibit inferior performance in all reasoning tasks. This suggests that while LLMs possess a broad knowledge base and strong language understanding,

Table 3: The deduction results of ProofWriter tasks (ChatGPT). Results are in %.

Category	Baseline	depth-1	depth-2	depth-3	depth-5
Symbols	Zero-Shot	69.1	62.3	59.4	52.8
	Zero-Shot-CoT	56.2	49.4	45.2	38.6
	Few-Shot-CoT	65.8	58.1	57.8	45.9
Semantics	Zero-Shot	69.0	63.5	60.3	51.4
	Zero-Shot-CoT	51.5	45.8	40.3	30.9
	Few-Shot-CoT	62.5	56.7	56.9	47.8

symbolic reasoning is not their primary strength compared to methods specifically designed for symbolic reasoning. This also suggests the potential of future neural-symbolic AI systems.

4.3 More Fine-grained Analysis about Semantics

306

307

308

309

310

311

312

317

318

319

320

321

322

323

324

325

326

331

332

333

334

335

338

339

340

342

345

346

347

348

349

350

352

353

354

355

356

The aforementioned experiments offer initial evidence highlighting the significance of semantics in the reasoning of LLMs. To further investigate this observation, we examine the influence of commonsense knowledge stored within LLMs on their semantic reasoning performance. Specifically, we explore three aspects: First, we examine the influence of commonsense knowledge stored within LLMs on their semantic reasoning performance. To achieve this, we remain the semantics (as semantics can encompass commonsense knowledge) and remove all given logical rules (in deduction) and facts (in induction). Please refer to Appendix A for prompts. This forces the LLMs to rely solely on their prior commonsense knowledge to infer the answers and allows us to assess the extent to which LLMs can leverage their internal knowledge to reason effectively without explicit in-context knowledge. Second, we retain the semantics of the datasets but introduce counter-commonsense logical rules. This requires LLMs to leverage in-context new knowledge and navigate the reasoning process by strictly adhering to the new information conflicting with the old knowledge. We implement it by shuffling relations as new relation labels to construct a new counter-commonsense dataset. For instance, we replace "motherOf" with "sisterOf", "parentOf" with "brotherOf", and "female" with "male". Consequently, for a rule such as $\forall x,y: \mathsf{parentOf}(x,y) \land \mathsf{female}(x) \to \mathsf{motherOf}(x,y)$, we obtain $\forall x, y$: brotherOf $(x, y) \land male(x) \rightarrow sisterOf(x, y)$. Thirdly, we use a subset of the ProofWriter OWA datasets for depths 0, 1, 2, 3 and 5, which contains synthetic facts and rules despite written in natural language but irrelevant to commonsense (see Figure 2). These investigations allow us to gain deeper insights into the effect of semantics on the reasoning capabilities of LLMs.

When semantics are consistent with commonsense As shown in Table 4, in the deductive reasoning experiment, *Removing rules/facts* achieves comparable results to *Semantics*; in the inductive reasoning experiment, *Removing rules/facts* outperforms *Symbols*, achieving 35.7% in GPT-4. These findings suggest that LLMs can perform deductive reasoning comparably by leveraging their stored commonsense knowledge without using the provided semantic knowledge, and providing symbolic instead of semantic knowledge in induction might even hurt the performance. Besides, GPT-4 significantly outperforms ChatGPT across all evaluation settings. The results may be attributed to the fact that the stored commonsense knowledge within GPT-4 is likely more extensive than that in ChatGPT or GPT-4 potentially possesses stronger reasoning capabilities. Additionally, there is a possibility of potential data contamination in the training process of GPT-4. For example, it has been trained on datasets, such as ProofWriter, which influenced the results.

When semantics are not consistent with commonsense To investigate the impact of semantics that are not consistent with commonsense, we introduce counter-commonsense (Counter-CS) scenarios, which is also shown in table 4. In comparison to Semantics and Symbols, we find that Counter-Commonsense performs worse than Semantics, even Symbols.

When semantics are not consistent with commonsense To investigate the impact of semantics that Table 4: Semantics, removing rules/facts and countercommonsense reasoning experiments (ChatGPT and GPT-4). Results are in %.

	deductive (F	ew-Shot-CoT)	inductive (Ze	ero-Shot-CoT)
	ChatGPT	GPT-4	ChatGPT	GPT-4
Semantics	71.8	90.0	25.0	53.6
Symbols	53.7	67.6	7.14	21.4
Remove R/F	70.1	90.4	7.14	35.7
Counter-CS	48.9	73.4	7.14	17.8

These findings suggest that when the in-context new knowledge conflicts with commonsense, LLMs struggle to accurately reason and predict.

When semantics are irrelevant to commonsense We use the ProofWriter tasks to test whether unmeaningful semantics are still useful. The results are shown in table 3. The *Symbols* setting performs comparably to the *Semantics* setting in the zero-shot setting, suggesting that when semantics are irrelevant to commonsense, they have little effect on the reasoning abilities of LLMs. In other words, when the task does not require deep semantic understanding or relies minimally on commonsense knowledge, the presence or absence of semantics does not significantly impact the performance of LLMs. However, in the CoT settings, we observe that *Semantics* is significantly worse than *Symbols*. This might be because step-by-step reasoning magnifies the disturbing effect brought by weird semantics such as "The squirrel needs the dog". Additionally, we observe that the CoT settings even perform worse than the zero-shot setting, with a higher frequency of the answer

"Cannot be determined.". Similar phenomenons are also observed in table 2, indicating that CoT may not be always helpful for reasoning tasks with in-context new knowledge.

More analysis and discussions

361

362

363

364

365

366

367

391

- (1) Induction and abduction underperform deduction: We compare the reasoning abilities of LLMs across induction and abduction tasks and find that they perform notably worse compared to deduction, regardless of whether semantics or symbols are used. When semantics are decoupled, the drop in performance is even more significant. These findings highlight the considerable room for improvement in LLMs' reasoning abilities and suggest that relying solely on semantics to achieve symbolic reasoning is challenging.
- (2) Shorter in-context knowledge enhances reasoning performance: To examine the influence 368 of context length on reasoning, we conducted an abductive reasoning experiment using a smaller 369 Symbolic Tree, containing approximately 12 entities and 100 facts. The results, provided in Ap-370 pendix P, show that abductive reasoning with a shorter context leads to better performance compared to a longer context. Besides, we also conduct deduction and induction experiments where LLMs are directly provided with the relevant facts related to the predicted fact or the predicted rule. The results are presented in Appendix K. This finding suggests that LLMs struggle with processing excessively long in-context information, particularly in reasoning tasks. The length of the context influences 375 376 reasoning performance, as shorter contexts make it easier to select relevant and useful information while minimizing the impact of unrelated content. 377
- (3) Effectiveness of commonsense expressed in natural language: We explore the representation 378 of knowledge in natural language and logic language forms in our experiments. The results, presented 379 in Appendix M, indicate that for tasks involving semantics, natural language descriptions are more 380 381 effective than logical language representations. Conversely, for symbolic and counter-commonsense 382 tasks, logic language performs better. This observation suggests that natural language representations better stimulate the semantic understanding capabilities of LLMs, while logical language 383 representations are more conducive to symbolic reasoning. 384
- (4) Zero-shot capabilities are approaching zero-shot-CoT capabilities: In Symbols setting, 385 comparing zero-shot with zero-shot-CoT across deduction, induction, and abduction evaluations, 386 we observe that zero-shot-CoT only marginally improves the performance compared to zero-shot 388 learning. This finding suggests that the zero-shot capabilities of current LLMs are approaching their zero-shot-CoT learning abilities. One plausible explanation is that ChatGPT has already been trained 389 on similar tasks with CoT and has memorized the instructions. Consequently, it implicitly follows 390 these instructions when applied to the same queries, even without explicit CoT guidance [76].
- (5) Utilizing internal knowledge outperforms external in-context knowledge: To explore the 392 393 ability of LLMs to utilize internal and external knowledge, we conduct an additional experiment where we provide LLMs with only the relevant facts related to the predicted fact. We compare the 394 performance of *Removing rules* (leveraging internal knowledge) with *Semantics* (providing external 395 logical rules). Surprisingly, we find that *Removing rules* performed better than *Semantics*. This 396 suggests that LLMs possess the necessary internal knowledge to support answering questions and 397 reasoning tasks, and leveraging this internal knowledge is more effective for reasoning than relying 398 on external logical rules. Detailed results and case studies can be found in Appendix K.1. 399

5 **Conclusion and Discussion**

Our paper presents the first comprehensive investigation of the role of semantics in LLMs' in-context 401 reasoning abilities by decoupling semantics from in-context prompts. Experimental results suggest 402 that: When semantics are consistent with commonsense, LLMs perform fairly well; when semantics 403 are decoupled or counter-commonsense, LLMs struggle to solve the reasoning tasks by leveraging 404 in-context new knowledge. These findings reveal the importance of semantics in LLMs' reasoning 405 abilities and inspire further research on unveiling the magic existing within the black-box LLMs. In 406 light of the findings identified in our analysis, we point out several potential future directions for the 407 development of large foundation models:

More complex symbolic reasoning benchmark: To improve LLMs' in-context symbolic reasoning abilities, developing new datasets with decoupled semantics and more complex reasoning tasks is necessary. These benchmarks should challenge LLMs with diverse and intricate symbolic knowledge.

Combination with external non-parametric knowledge base: As our experimental results show, the memorization abilities of LLMs are not comparable to existing graph-based methods. This motivates integrating LLMs with external non-parametric knowledge bases, such as graph databases, to enhance their knowledge insertion and updating. This hybrid approach can leverage the strengths of LLMs' language understanding and the comprehensive, accurate and up-to-date knowledge stored in non-parametric sources.

Improving the ability of processing in-context knowledge: We observe that LLMs perform better under shorter context (discussion 4.5 (2)) and when only provided with the relevant facts related to the question (discussion 4.5 (6)). The results indicate that LLMs probably struggle with processing excessively long in-context information. Moreover, discussion 4.5 (5) also suggests LLMs are more reliable to leverage internal knowledge. As a result, effectively utilizing external (in-context) knowledge to perform situated tasks remains an important challenge for LLMs. This includes developing mechanisms to better encode and retrieve relevant information from the in-context knowledge.

426 References

- [1] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
 Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners.
 Advances in neural information processing systems, 33:1877–1901, 2020.
- 430 [2] Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in large language models. *arXiv preprint arXiv:2210.03493*, 2022.
- 432 [3] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S
 433 Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of
 434 foundation models. arXiv preprint arXiv:2108.07258, 2021.
- 435 [4] Kevin Lu, Aditya Grover, Pieter Abbeel, and Igor Mordatch. Pretrained transformers as universal computa-436 tion engines. *arXiv preprint arXiv:2103.05247*, 1, 2021.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. Chain
 of thought prompting elicits reasoning in large language models. arXiv preprint arXiv:2201.11903, 2022.
- 439 [6] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *arXiv preprint arXiv:2205.11916*, 2022.
- 441 [7] Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke
 442 Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? *arXiv preprint*443 *arXiv:2202.12837*, 2022.
- 444 [8] Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. What makes good in-context examples for gpt-3? *arXiv preprint arXiv:2101.06804*, 2021.
- [9] Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. Calibrate before use: Improving few-shot performance of language models. In *International Conference on Machine Learning*, pages 12697–12706. PMLR, 2021.
- 449 [10] Timo Schick and Hinrich Schütze. Exploiting cloze questions for few shot text classification and natural language inference. *arXiv preprint arXiv:2001.07676*, 2020.
- 451 [11] Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. Fantastically or-452 dered prompts and where to find them: Overcoming few-shot prompt order sensitivity. *arXiv preprint* 453 *arXiv:2104.08786*, 2021.
- 454 [12] Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661*, 2022.
- 456 [13] Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context learning as implicit bayesian inference. *arXiv preprint arXiv:2111.02080*, 2021.
- 458 [14] Steven A Sloman. The empirical case for two systems of reasoning. *Psychological bulletin*, 119(1):3, 1996.
- [15] Thomas Gilovich, Dale Griffin, and Daniel Kahneman. Heuristics and biases: The psychology of intuitive
 judgment. Cambridge university press, 2002.
- [16] Oyvind Tafjord, Bhavana Dalvi Mishra, and Peter Clark. Proofwriter: Generating implications, proofs, and
 abductive statements over natural language. arXiv preprint arXiv:2012.13048, 2020.
- Robin Cooper. The fracas consortium robin cooper, dick crouch, jan van eijck, chris fox, josef van genabith, jan jaspars, hans kamp, david milward, manfred pinkal, massimo poesio, steve pulman. 1996.
- 466 [18] SamuelR. Bowman, Gabor Angeli, Christopher Potts, and ChristopherD. Manning. A large annotated corpus for learning natural language inference, Aug 2015.
- 468 [19] Adina Williams, Nikita Nangia, and SamuelR. Bowman. A broad-coverage challenge corpus for sentence understanding through inference, Apr 2017.
- 470 [20] Jean-Philippe Bernardy and Stergios Chatzikyriakidis. Improving the precision of natural textual entailment problem datasets., May 2020.
- 472 [21] Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. Commonsenseqa: A question answering challenge targeting commonsense knowledge, Nov 2018.

- Shourya Aggarwal, Divyanshu Mandowara, Vishwajeet Agrawal, Dinesh Khandelwal, Parag Singla, and Dinesh Garg. Explanations for commonsenseqa: New dataset and models. In *Proceedings of the S9th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, Jul 2021. doi: 10.18653/v1/2021.
 acl-long.238. URL http://dx.doi.org/10.18653/v1/2021.acl-long.238.
- 479 [23] Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. Drop: 480 A reading comprehension benchmark requiring discrete reasoning over paragraphs, Mar 2019.
- 481 [24] Wenhu Chen, Hanwen Zha, Zhiyu Chen, Wenhan Xiong, Hong Wang, and William Yang Wang. Hybridqa:
 482 A dataset of multi-hop question answering over tabular and textual data, Apr 2020.
- Zonglin Yang, Li Dong, Xinya Du, Hao Cheng, Erik Cambria, Xiaodong Liu, Jianfeng Gao, and Furu Wei.
 Language models as inductive reasoners. arXiv preprint arXiv:2212.10923, 2022.
- 485 [26] Abulhair Saparov and He He. Language models are greedy reasoners: A systematic formal analysis of chain-of-thought. *arXiv* preprint arXiv:2210.01240, 2022.
- [27] Simeng Han, Hailey Schoelkopf, Yilun Zhao, Zhenting Qi, Martin Riddell, Luke Benson, Lucy Sun,
 Ekaterina Zubova, Yujie Qiao, Matthew Burtell, et al. Folio: Natural language reasoning with first-order
 logic. arXiv preprint arXiv:2209.00840, 2022.
- 490 [28] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- 492 [29] Antonia Creswell, Murray Shanahan, and Irina Higgins. Selection-inference: Exploiting large language
 493 models for interpretable logical reasoning. arXiv preprint arXiv:2205.09712, 2022.
- 494 [30] Subhro Roy and Dan Roth. Solving general arithmetic word problems. *arXiv preprint arXiv:1608.01413*, 495 2016.
- [31] Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. Did aristotle
 use a laptop? a question answering benchmark with implicit reasoning strategies. *Transactions of the Association for Computational Linguistics*, 9:346–361, 2021.
- [32] Ishita Dasgupta, Andrew K Lampinen, Stephanie CY Chan, Antonia Creswell, Dharshan Kumaran, James L
 McClelland, and Felix Hill. Language models show human-like content effects on reasoning. arXiv preprint
 arXiv:2207.07051, 2022.
- 502 [33] Viktor Schlegel, Kamen V Pavlov, and Ian Pratt-Hartmann. Can transformers reason in fragments of natural language? *arXiv preprint arXiv:2211.05417*, 2022.
- Rajarshi Das, Manzil Zaheer, Dung Thai, Ameya Godbole, Ethan Perez, Jay Yoon Lee, Lizhen Tan,
 Lazaros Polymenakos, and Andrew McCallum. Case-based reasoning for natural language queries over
 knowledge bases. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 9594–9611, Online and Punta Cana, Dominican Republic, November 2021. Association
 for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.755. URL https://aclanthology.org/2021.emnlp-main.755.
- 510 [35] Ohad Rubin, Jonathan Herzig, and Jonathan Berant. Learning to retrieve prompts for in-context learning.
 511 arXiv preprint arXiv:2112.08633, 2021.
- [36] Hongjin Su, Jungo Kasai, Chen Henry Wu, Weijia Shi, Tianlu Wang, Jiayi Xin, Rui Zhang, Mari Ostendorf,
 Luke Zettlemoyer, Noah A. Smith, and Tao Yu. Selective annotation makes language models better
 few-shot learners, 2022.
- [37] Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang Sutawika, Zaid Alyafeai, Antoine 515 Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish 516 Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal Nayak, Deba-517 jyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin 518 Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, 519 Andrea Santilli, Thibault Fevry, Jason Alan Fries, Ryan Teehan, Tali Bers, Stella Biderman, Leo Gao, 520 Thomas Wolf, and Alexander M. Rush. Multitask prompted training enables zero-shot task generalization, 521 522 2022.
- [38] Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M.
 Dai, and Quoc V. Le. Finetuned language models are zero-shot learners, 2022.

- Yao Lu, Max Bartolo, AlastairP. Moore, Sebastian Riedel, and Pontus Stenetorp. Fantastically ordered
 prompts and where to find them: Overcoming few-shot prompt order sensitivity, Apr 2021.
- 527 [40] Albert Webson and Ellie Pavlick. Do prompt-based models really understand the meaning of their 528 prompts? In Proceedings of the 2022 Conference of the North American Chapter of the Association for 529 Computational Linguistics: Human Language Technologies, pages 2300–2344, Seattle, United States, 530 July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.167. URL 531 https://aclanthology.org/2022.naacl-main.167.
- [41] Boshi Wang, Sewon Min, Xiang Deng, Jiaming Shen, You Wu, Luke Zettlemoyer, and Huan Sun.
 Towards understanding chain-of-thought prompting: An empirical study of what matters. arXiv preprint
 arXiv:2212.10001, 2022.
- Miles Turpin, Julian Michael, Ethan Perez, and Samuel R Bowman. Language models don't always say
 what they think: Unfaithful explanations in chain-of-thought prompting. arXiv preprint arXiv:2305.04388,
 2023.
- 538 [43] Aman Madaan and Amir Yazdanbakhsh. Text and patterns: For effective chain of thought, it takes two to tango. *arXiv preprint arXiv:2209.07686*, 2022.
- 540 [44] Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke 541 Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work?
- [45] Junyeob Kim, Hyuhng Joon Kim, Hyunsoo Cho, Hwiyeol Jo, Sang-Woo Lee, Sang-goo Lee, Kang Min
 Yoo, and Taeuk Kim. Ground-truth labels matter: A deeper look into input-label demonstrations. arXiv
 preprint arXiv:2205.12685, 2022.
- [46] Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu,
 Da Huang, Denny Zhou, et al. Larger language models do in-context learning differently. arXiv preprint
 arXiv:2303.03846, 2023.
- Jerry Wei, Le Hou, Andrew Lampinen, Xiangning Chen, Da Huang, Yi Tay, Xinyun Chen, Yifeng Lu,
 Denny Zhou, Tengyu Ma, et al. Symbol tuning improves in-context learning in language models. arXiv
 preprint arXiv:2305.08298, 2023.
- Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models, Nov 2022.
- 553 [49] George Boole. The mathematical analysis of logic. Philosophical Library, 1847.
- [50] John McCarthy. Recursive functions of symbolic expressions and their computation by machine, part i.
 Communications of the ACM, 3(4):184–195, 1960.
- [51] Norbert Fuhr. Probabilistic datalog: Implementing logical information retrieval for advanced applications.
 Journal of the American Society for Information Science, 51(2):95–110, 2000.
- 558 [52] Thomas Eiter, Giovambattista Ianni, and Thomas Krennwallner. Answer set programming: A primer.559 Springer, 2009.
- 560 [53] Kexin Yi, Jiajun Wu, Chuang Gan, Antonio Torralba, Pushmeet Kohli, and Josh Tenenbaum. Neural-561 symbolic vqa: Disentangling reasoning from vision and language understanding. *Advances in neural* 562 *information processing systems*, 31, 2018.
- John McCarthy and Patrick J Hayes. Some philosophical problems from the standpoint of artificial
 intelligence. In *Readings in artificial intelligence*, pages 431–450. Elsevier, 1981.
- 565 [55] Nada Lavrac and Saso Dzeroski. Inductive logic programming. In WLP, pages 146–160. Springer, 1994.
- [56] Allen Newell and Herbert A Simon. Computer science as empirical inquiry: Symbols and search. In ACM
 Turing award lectures, page 1975. 2007.
- [57] David Landy, Colin Allen, and Carlos Zednik. A perceptual account of symbolic reasoning. Frontiers in
 psychology, 5:275, 2014.
- 570 [58] Megan Fowler. A human-centric system for symbolic reasoning about code. 2021.
- 571 [59] Philip N Johnson-Laird. Deductive reasoning. Annual review of psychology, 50(1):109–135, 1999.
- 572 [60] Gyöngyi Kovács and Karen M Spens. Abductive reasoning in logistics research. *International journal of physical distribution & logistics management*, 2005.

- 574 [61] Jing Qian, Hong Wang, Zekun Li, Shiyang Li, and Xifeng Yan. Limitations of language models in arithmetic and symbolic induction. *arXiv preprint arXiv:2208.05051*, 2022.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch,
 Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. Beyond the imitation game:
 Quantifying and extrapolating the capabilities of language models. arXiv preprint arXiv:2206.04615, 2022.
- Eui Chul Shin, Illia Polosukhin, and Dawn Song. Improving neural program synthesis with inferred execution traces. *Advances in Neural Information Processing Systems*, 31, 2018.
- 581 [64] Guillaume Lample and François Charton. Deep learning for symbolic mathematics. *arXiv preprint* arXiv:1912.01412, 2019.
- Vishal Pallagani, Bharath Muppasani, Keerthiram Murugesan, Francesca Rossi, Lior Horesh, Biplav
 Srivastava, Francesco Fabiano, and Andrea Loreggia. Plansformer: Generating symbolic plans using
 transformers. arXiv preprint arXiv:2212.08681, 2022.
- [66] Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham
 Neubig. Pal: Program-aided language models. In *International Conference on Machine Learning*, pages
 10764–10799. PMLR, 2023.
- Feter Cathcart Wason and Philip Nicholas Johnson-Laird. Psychology of reasoning: Structure and content,
 volume 86. Harvard University Press, 1972.
- [68] P. C. Wason. Reasoning about a rule. *Quarterly Journal of Experimental Psychology*, page 273–281, Jul
 2007. doi: 10.1080/14640746808400161. URL http://dx.doi.org/10.1080/14640746808400161.
- Fig. [69] Ronald Fagin, Joseph Y Halpern, Yoram Moses, and Moshe Vardi. Reasoning about knowledge. MIT press, 2004.
- 595 [70] Douglas N Walton. What is reasoning? what is an argument? *The journal of Philosophy*, 87(8):399–419, 1990.
- [71] Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. *Advances in neural information processing systems*, 26, 2013.
- 600 [72] Patrick Hohenecker and Thomas Lukasiewicz. Ontology reasoning with deep neural networks. *Journal of Artificial Intelligence Research*, 68:503–540, 2020.
- 602 [73] Massimo Carro. Nosql databases. arXiv preprint arXiv:1401.2101, 2014.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford alpaca, 2023.
 - [75] Luis Galárraga, Christina Teflioudi, Katja Hose, and Fabian M Suchanek. Fast rule mining in ontological knowledge bases with amie

+

- +. The VLDB Journal, 24(6):707–730, 2015.
- [76] Jiuhai Chen, Lichang Chen, Heng Huang, and Tianyi Zhou. When do you need chain-of-thought prompting for chatgpt? *arXiv preprint arXiv:2304.03262*, 2023.

Table of Contents

610	A	Pron	npts		16
611		A.1	Deducti	ive reasoning	16
612			A.1.1	Zero-Shot	16
613			A.1.2	Zero-Shot-CoT	16
614			A.1.3	Few-Shot-CoT	17
615			A.1.4	Examples of Few-Shot-CoT Demonstrations	17
616			A.1.5	Zero-Shot of removing rules setting	17
617		A.2	Inductiv	ve reasoning	18
618			A.2.1	Zero-Shot	18
619			A.2.2	Zero-Shot CoT	18
620			A.2.3	Zero-Shot of removing facts setting	18
621		A.3	Abducti	ive reasoning	18
622			A.3.1	Zero-Shot	18
623			A.3.2	Zero-Shot-CoT	19
624			A.3.3	Few-Shot-CoT prompt	19
625			A.3.4	Examples of Few-Shot-CoT Demonstrations	19
626			A.3.5	Few-Shot CoT of removing rules setting	20
627		A.4	Commo	onsense reasoning	21
628			A.4.1	Zero-Shot prompt	21
629			A.4.2	Few-Shot prompt	21
630				Examples of Few-Shot Demonstrations	21
631		A.5		A Fine-tuning Prompt	21
632	B	Dedu	iction ex	xamples of Symbolic Tree datasets	21
633		B.1		ics	22
634			B.1.1	Logic language representations	22
635			B.1.2	Natural language representations	23
636		B.2	Symbol	lization	25
637			B.2.1	Logic language representations	25
638			B.2.2	Natural language representations:	27
639		B.3	Semant	ics of removing rule setting	28
640	C	Exan	nples of	ProofWriter	29
641		C.1	Semant	ics	29
642		C.2	Symbol	ls	30
643				ro-Shot prompting of memorization abilities of neural-based and symbolic-based methods	31

645	F	Introduction of Neo4j	32
646	G	Task definitions	32
647	Н	Implementation of memorization	33
648	I	Consistency of knowledge base	33
649	J	Different Few-Shot-CoT examples	33
650		J.1 Few-Shot-CoT examples with Auto-Filter	33
651		J.2 Few-Shot-CoT examples with Manually-Designed	34
652	K	Reasoning results of after human selection	36
653		K.1 Comparison of Removing rules and Semantics	36
654		K.2 Reasoning results after human selection	36
655	L	Reasoning with natural language	37
656	M	Reasoning results of two representations	37
657	N	Reasoning with garbled symbols	37
658	o	More reasoning results of Symbolic Tree	39
659	P	Abductive reasoning on smaller datasets	39
660	Q	Replacing entity labels	40
661	R	Multi-short rules	40
662	A	Prompts	
663	Α.	1 Deductive reasoning	
664	A.	1.1 Zero-Shot	
665		stem: You are a helpful assistant with deductive reasoning abilities.	
666 667	use	er: I will provide a set of logical rules L1 to L{number of rules} and facts F1 to F{numb of basic facts}. Please select one single logical rule from L1 to L{number of rules} a	
668		a few facts from F1 to F{number of basic facts} to predict True/False of the unknown f	
669		using deductive reasoning.	
670		gical_rules: {logical rules}	
671		cts: {basic facts}	
672 673		known fact: {statement} e answer (True or False) is:	
674	Α.	1.2 Zero-Shot-CoT	
675	sys	stem: You are a helpful assistant with deductive reasoning abilities. Please select one	
676 677	บรร	single logical rule and a few facts to predict True/False of the following statement. er: I will provide a set of logical rules L1 to L{number of rules} and facts F1 to F{number.	er
678	200	of basic facts}. Please select one single logical rule from L1 to L{number of rules} a	
679		a few facts from F1 to F{number of basic facts} to predict True/False of the following	
680 681	Log	statement using deductive reasoning. gical rules: {logical rules}	
	_		

```
Facts: {basic facts}
682
683
    Statement: {statement}
    Answer with True or False? Let's think step by step.
```

A.1.3 Few-Shot-CoT

685

```
system: You are a helpful assistant with deductive reasoning abilities. Please select one
686
         single logical rule and a few facts to predict True/False of the following statement.
687
     user: I will provide a set of logical rules L1 to L{number of rules} and facts F1 to F{number
688
          of basic facts}
     Logical rules: {logical rules}
690
     Facts: {basic facts}
691
     Please select one single logical rule from L1 to L{number of rules} and a few facts from F1 to
692
          F{number of basic facts} to predict True/False of the following statement using
693
     {Demonstrations}
695
     Statement: {statement}
696
     Answer:
697
```

Examples of Few-Shot-CoT Demonstrations A.1.4

```
Statement: r14(Amelie, Jonathan)
699
     Answer: We can use L11: \forall A, B, C, D, E : r3(B, A) \land r3(B, C) \land r3(C, D) \land r3(D, E) \land r2(A) \rightarrow
700
     r14(A, E) to make a deduction. From the facts, we know that r3 (Amelie, Thomas) and r3 (Thomas,
701
     Jonathan). Therefore, we can apply L11 with A = Amelie, B = Thomas, C = Jonathan, D and E as
702
     variables. This gives us: r3(Thomas, Jonathan) \wedge r3(Amelie, Thomas) \wedge r3(Thomas, D) \wedge r3(D, E)
703
     \wedge r2(Amelie) \rightarrow r14(Amelie, E) Since we know that r2(Amelie) is true from F7, we can simplify
704
     the above statement to: r3(Thomas, Jonathan) \wedge r3(Amelie, Thomas) \wedge r3(Thomas, D) \wedge r3(D, E)
705
     \rightarrow r14(Amelie, E) Now, we can substitute the facts r3(Amelie, Thomas) and r3(Thomas, Jonathan)
706
     to get: r3(Thomas, Jonathan) \wedge r3(Thomas, D) \wedge r3(D, E) \rightarrow r14(Amelie, E) We can see that this
707
     statement is true if we choose D = Alina and E = Jonathan, since we know that r_3 (Thomas, Alina)
708
     from F50 and r3(Alina, Jonathan) from F56. Therefore, the statement r14(Amelie, Jonathan) is true.
709
```

- Statement: r31(Philipp, Nina) 710
- Answer: Let's use L28: $\forall A, B, C, D : r3(B, A) \land r3(C, B) \land r3(C, D) \land r1(A) \rightarrow r31(A, D)$ to 711 see if we can prove the statement r31(Philipp, Nina). From the facts, we know that r3(Luisa, Nina)(F52) and r3(Emilia, Philipp) (F62). We can use L28 with A = Philipp, B = Emilia, C = Luisa, and D = Nina to get: $r3(\text{Emilia}, \text{Philipp}) \land r3(\text{Luisa}, \text{Emilia}) \land r3(\text{Luisa}, \text{Nina}) \land r1(\text{Philipp}) \rightarrow$ r31(Philipp, Nina) Since all the conditions are met, we can conclude that r31(Philipp, Nina) is True. 715 Therefore, the answer is True. 716
- Statement: r8(Leonie, Nico) 717
- Answer: We can use L5 to make a deduction about the statement r8 (Leonie, Nico). L5 states that if 718 there exist three individuals A, B, and C such that A is related to B, B is related to C, and A has a certain property (r2), then A is related to C in a certain way (r8). Using this rule, we can see that we 720 have the following facts: -r3(Leonie, Emily) 721 - r3(Emily, Nico) 722
- Therefore, we can conclude that r8(Leonie, Nico) is true. 723
- Statement: r7(Patrick, Alina) 724
- Answer: We can use L4: $\forall A, B : r3(A, B) \land r1(A) \rightarrow r7(A, B)$ to determine if r7(Patrick, Alina) 725 is true or false. From the facts, we know that r^3 (Patrick, Alina) is true (F55). We also know that r1(Patrick) is true (F9). Therefore, we can apply L4 to conclude that r7(Patrick, Alina) is true. 727
- Answer: True. 728
- Statement: r27(Jonathan, Leonie) 729
- Answer: Let's use L4: $\forall A, B : r3(A, B) \land r1(A) \rightarrow r7(A, B)$ and F56: r3(Patrick, Jonathan) to infer 730
- that r7(Patrick, Jonathan) is true. Then, we can use L24: $\forall A, B, C : r3(B, A) \land r3(C, B) \land r1(A) \rightarrow$ 731
- r27(A, C) and F44: r3(Leonie, Emily) to infer that r27(Jonathan, Leonie) is false, since there is no 732
- fact that supports r3(Jonathan, Emily). Therefore, the answer is False.

A.1.5 Zero-Shot of removing rules setting

```
735    system: Please answer the question only with True or False.
736    user: I will provide a set of facts. Please predict True/False of the unknown fact based on
737    given facts.
738    Facts: {facts}
739    Unknown fact: {statement}
740    The answer (True or False) is:
```

741 A.2 Inductive reasoning

742 A.2.1 Zero-Shot

```
system: You are a helpful assistant with inductive reasoning abilities. Please generate one
743
         single rule to match the template and logically entail the facts. Note that the symbol
744
         '##' in the template should be filled with either 'r1' or 'r45', while the symbol '++'
745
         should be filled with either 'r43' or 'r44'.
746
    user: I will give you a set of facts F1 to F{number of basic facts}, facts G1 to G{number of
747
748
         inferred fact} and a template for a logical rule. Please generate one single rule to
749
         match the template and logically entail the facts G1 to G{number of inferred fact} based
         on facts F1 to F{number of basic facts}.
750
    Facts: {facts}
751
752
    Template: {rule template}
    Note that the symbol '##' in the template should be filled with either 'r1' or 'r45', while
753
         the symbol '++' should be filled with either 'r43' or 'r44'.
754
    After filling in the template, the generated rule is:
755
```

756 A.2.2 Zero-Shot CoT

```
system: You are a helpful assistant with inductive reasoning abilities. Please generate one
757
758
          single rule to match the template and logically entail the facts. Note that the symbol
          '##' in the template should be filled with either 'r1' or 'r45', while the symbol '++'
759
760
          should be filled with either 'r43' or 'r44'.
    user: I will give you a set of facts F1 to F{number of basic facts}, facts G1 to G{number of
761
          inferred fact) and a template for a logical rule. Please generate one single rule to
762
         match the template and logically entail the facts G1 to G{number of inferred fact} based
763
         on facts F1 to F{number of basic facts}.
764
765
    Facts: {facts}
    Template: {rule template}
766
    Note that the symbol '##' in the template should be filled with either 'r1' or 'r45', while
767
768
          the symbol '++' should be filled with either 'r43' or 'r44'.
    After filling in the template, the generated rule is: Let's think step by step.
769
```

770 A.2.3 Zero-Shot of removing facts setting

```
system: Please generate one single rule to match the template. Note that the symbol '##' in
771
         the template should be filled with either 'parent' or 'child', while the symbol '++'
772
         should be filled with either 'male' or 'female'.
773
     user: I will give you a template for a logical rule. Please generate one single rule to match
          the template and logically infer the relation sister
775
     Template: If A is ## of B and B is ## of C and A is ++, then A is sister of C.
776
    Note that the symbol '##' in the template should be filled with either 'parent' or 'child',
777
         while the symbol '++' should be filled with either 'male' or 'female'.
778
     After filling in the template, the generated rule is:
```

780 A.3 Abductive reasoning

781 A.3.1 Zero-Shot

```
system: You are a helpful assistant with abductive reasoning abilities. Please select one
782
         single logical rule and a few facts to explain the following statement.
783
    user: I will provide a set of logical rules L1 to L{number of rules} and facts F1 to F{number
784
785
         of basic facts}. Please select one single logical rule from L1 to L{number of rules} and
786
         a few facts from F1 to F{number of basic facts} to explain the following statement.
    Rules: {logical rules}
787
    Facts: {basic facts}
788
    Statement: {statement}
789
    Answer with the numbers of the selected rule and facts. The selected rule and facts are:
```

A.3.2 Zero-Shot-CoT

```
system: You are a helpful assistant with abductive reasoning abilities. Please select one
792
         single logical rule and a few facts to explain the following statement.
793
794
    user: I will provide a set of logical rules L1 to L{number of rules} and facts F1 to F{number
795
         of basic facts}. Please select one single logical rule from L1 to L{number of rules} and
796
         a few facts from F1 to F{number of basic facts} to explain the following statement.
    Rules: {logical rules}
797
    Facts: {basic facts}
798
    Statement: {statement}
799
    Answer with the numbers of the selected rule and facts. The selected rule and facts are: Let's
800
801
          think step by step.
```

802 A.3.3 Few-Shot-CoT prompt

```
system: You are a helpful assistant with abductive reasoning abilities. Please select one
803
804
         single logical rule and a few facts to explain the following statement.
805
     user:
        I will provide a set of logical rules L1 to L{number of rules} and facts F1 to F{number of
806
807
             basic facts}. Please select one single logical rule from L1 to L{number of rules} and
             a few facts from F1 to F{number of basic facts} to explain the following statement.
808
809
        Rules: {logical rules}
        Facts: {basic facts}
810
811
        {Demonstrations}
812
        Statement: {statement}
        Answer:
813
```

814 A.3.4 Examples of Few-Shot-CoT Demonstrations

- 815 Statement: r4(Laura, Fabian)
- Answer: To explain the statement r4(Laura, Fabian), we can use the rule L3 and fact F1 and F37.
- 817 L3: $\forall A, B : r1(A, B) \land r44(A) \rightarrow r4(A, B)$
- 818 F1: r1(Laura, Fabian)
- 819 F37: r44(Laura)
- According to L3, if there is a relationship r1 between A and B, and A has the property r44, then the
- 821 relationship r4 holds between A and B. In this case, we have a relationship r1 between Laura and
- Fabian (F1) and Laura has the property r44 (F37). Therefore, the relationship r4 holds between Laura
- and Fabian, which is the given statement.
- The selected rule and facts are: L3, F1, and F37.
- Statement: r5(Stefan, Emily)
- Answer: To explain the statement r5(Stefan, Emily), we can use the rule L4 and fact F17. L4:
- 827 $\forall A, B : r1(A, B) \land r43(A) \rightarrow r5(A, B)$
- 828 F17: r1(Stefan,Emily)
- According to rule L4, if there is a relation r1(A, B) and A has the property r43, then the relation
- r5(A, B) holds. In our case, we have the fact F17, which states that there is a relation r1(Stefan,
- Emily). Additionally, Stefan has the property r43 (F61). Therefore, we can conclude that the relation
- 832 r5(Stefan, Emily) holds.
- The selected rule and facts are: L4, F17, and F61.
- Statement: r23(Felix, Elias)
- Answer: To explain the statement r23 (Felix, Elias), we can use the rule L22 and fact F5. L22:
- 836 $\forall A, B : r1(B, A) \land r43(A) \to r23(A, B)$
- F5: r1(Elias, Felix)
- According to rule L22, if there exists a relationship r1(B, A) and A has the property r43, then the
- relationship r23(A, B) holds. In our case, we have the fact F5, which states that there is a relationship
- r1(Elias, Felix). Additionally, we know that Felix has the property r43 (F59). Therefore, applying
- rule L22, we can conclude that the relationship r23 (Felix, Elias) holds.
- The selected rule and facts are: L22, F5, and F59.
- Statement: r3(Felix, Fabian)
- Answer: To explain the statement r3(Felix, Fabian), we can use the rule L2 and facts F2, F1, and F59.
- 845 L2: $\forall A, B, C : r1(B, A) \land r1(B, C) \land r43(A) \rightarrow r3(A, C)$
- F1: r1(Laura, Fabian)

```
F2: r1(Laura, Felix)
847
    F59: r43(Felix)
848
     By substituting the values into L2, we get:
849
    r1(Laura, Felix) \land r1(Laura, Fabian) \land r43(Felix) \rightarrow r3(Felix, Fabian)
    Since all the premises are true according to the given facts, the conclusion r3 (Felix, Fabian) is also
851
852
    The selected rule and facts are: L2, F1, F2, and F59.
     Statement: r7(Fabian, Helga)
854
    Answer: To explain the statement r7(Fabian, Helga), we can use the rule L6 and facts F13, F21, and
855
856
    L6: \forall A, B, C : r1(A, B) \land r1(B, C) \land r43(A) \rightarrow r7(A, C)
857
    F13: r1(Fabian,Thomas)
858
    F21: r1(Thomas,Helga)
859
    F39: r43(Fabian)
    According to L6, if there exist A, B, and C such that r1(A, B), r1(B, C), and r43(A), then r7(A, C)
     holds. In our case, A = Fabian, B = Thomas, and C = Helga. Since we have r1(Fabian, Thomas)
862
     from F13, r1 (Thomas, Helga) from F21, and r43 (Fabian) from F39, we can conclude that r7 (Fabian,
863
     Helga) is true.
864
    So, the selected rule and facts are: L6, F13, F21, and F39.
865
```

866 A.3.5 Few-Shot CoT of removing rules setting

```
system: You are a helpful assistant with abductive reasoning abilities. Please select a few
867
868
         facts to explain the following statement.
    user: I will provide a set of facts F1 to F63. Please select a few facts from F1 to F63 to
869
870
         explain the following statement.
871
    Facts: {facts}
872
873
    Statement: Laura is mother of Felix.
    Answer: To explain the statement "Laura is mother of Felix", we can use Facts:
874
    Fact F2 states: Laura is parent of Felix.
875
    Fact F37 states: Laura is female.
    Using F2 and F37, we can conclude that "Laura is mother of Felix" holds.
877
    Therefore, the selected rule and facts are F2, F37.
878
879
880
    Statement: Samuel is brother of Alina.
881
    Answer: To infer the statement "Samuel is brother of Alina", we have:
    F27: Patrick is parent of Samuel.
882
    F28: Patrick is parent of Alina.
884
    F47: Samuel is male.
    Based on these facts, we can infer "Samuel is brother of Alina":
885
    Therefore, the selected rule and facts are F27, F28, F47.
887
888
    Statement: Patrick is grandfather of David.
889
    Answer: To explain the statement "Patrick is grandfather of David", we have:
    F28: Patrick is parent of Alina.
890
    F7: Alina is parent of David.
    F45: Patrick is male.
892
    Based on these facts, we can infer "Patrick is grandfather of David":
893
894
    Therefore, the selected rule and facts are F28, F7, F45.
895
    Statement: Amelie is daughter of Elena.
    Answer: To explain the statement "Amelie is daughter of Elena", we have:
897
    F20: Elena is parent of Amelie.
898
899
    F43: Amelie is female.
    Based on these facts, we can infer "Amelie is daughter of Elena".
900
    Therefore, the selected rule and facts are F20, F43.
902
    Statement: Claudia is sister of Felix
903
904
    Answer: To prove the statement "Claudia is sister of Felix", we can use facts:
    F3: Laura is parent of Claudia.
905
    F2: Laura is parent of Felix.
907
    F40: Claudia is female.
    Based on these facts, we can infer "Claudia is sister of Felix".
908
    Therefore, the selected rule and facts are F3, F2, F40.
910
    Statement: Laura is mother of Fabian.
```

912 Answer:

913 A.4 Commonsense reasoning

914 A.4.1 Zero-Shot prompt

```
915     system: Please choose answer from [A, B, C, D, E] and output one capital letter only.
916     user: Question: {Question}
917     {Multiple choices}
918     Answer:
```

919 A.4.2 Few-Shot prompt

```
920     system: Please choose answer from [A, B, C, D, E] and output one capital letter only.
921     user: {Demonstrations:}
922     Question: {Question}
923     {Multiple choices}
924     Answer:
```

925 A.4.3 Examples of Few-Shot Demonstrations

926 Standard:

- 927 Question: When is hotter?
- 928 Choices: A.desert B.night C.morning D.noon
- 929 Answer: A.

930 Predict-explain:

- 931 Question: When is hotter?
- 932 Choices: A.desert B.night C.morning D.noon
- Answer: A. The desert is hotter than the other options listed because it receives direct sunlight and
- has little to no vegetation to provide shade or cool the air. During the day, the sun beats down on the
- desert, causing temperatures to rise significantly. At night, the lack of vegetation and moisture in the
- 936 air means that heat is quickly lost, causing temperatures to drop rapidly. In the morning and at noon,
- the sun is still high in the sky, but temperatures may be slightly cooler due to the presence of some
- shade and cooler air from the night before.

939 Expain-predict:

948

955

- 940 Question: When is hotter?
- 941 Choices: A.desert B.night C.morning D.noon
- The desert is hotter than the other options listed because it receives direct sunlight and has little to no
- vegetation to provide shade or cool the air. During the day, the sun beats down on the desert, causing
- temperatures to rise significantly. At night, the lack of vegetation and moisture in the air means that
- heat is quickly lost, causing temperatures to drop rapidly. In the morning and at noon, the sun is still
- 946 high in the sky, but temperatures may be slightly cooler due to the presence of some shade and cooler
- 947 air from the night before. Answer: A.

A.5 LLaMA Fine-tuning Prompt

```
949 Below is an instruction that describes a task, paired with an input that provides further
950 context.
951 Write a response that appropriately completes the request.
952 Instruction: {Head} is the {Relation} of {Tail}
953 Input: {input}
954 Response:
```

B Deduction examples of Symbolic Tree datasets

- 956 In this section, we provide examples of deduction experiments conducted on the Symbolic Tree
- 957 datasets. We present examples for both the Semantics and Symbols settings, represented in both
- 958 natural language text and logic language

59 B.1 Semantics

960 B.1.1 Logic language representations

```
Logical rules:
 961
         L1: $\forall A,B,C: parentOf(B, A) \land parentOf(B, C) \land female(A) \rightarrow sisterOf(A,
 962
 963
         L2: $\forall A,B,C: parentOf(B, A) \land parentOf(B, C) \land male(A) \rightarrow brotherOf(A,
 964
 965
                 C)$
         L3: $\forall A,B: parentOf(A, B) \land female(A) \rightarrow motherOf(A,B)$
 966
         L4: $\forall A,B: parentOf(A, B) \land male(A) \rightarrow fatherOf(A,B)$
 967
 968
         L5: $\forall A,B,C: parentOf(A, B) \land parentOf(B, C) \land female(A) \rightarrow
                 grandmotherOf(A,C)$
 969
 970
         L6: $\forall A,B,C: parentOf(A, B) \land parentOf(B, C) \land male(A) \rightarrow
                 grandfatherOf(A,C)$
 971
         L7: $\forall A,B,C,D: parentOf(A, B) \land parentOf(B, C) \land parentOf(C, D) \land female(A)
 972
 973
                  \rightarrow greatGrandmotherOf(A,D)$
         L8: f(A, B) \rightarrow f(B, C) \rightarrow f(C, D) \rightarrow f(A, B)
 974
 975
                 rightarrow greatGrandfatherOf(A,D)$
         L9: $\forall A,B,C,D: parentOf(B, A) \land parentOf(B, C) \land parentOf(C, D) \land female(A)
 976
 977
                  \rightarrow auntOf(A,D)$
 978
         L10: $\forall A,B,C,D: parentOf(B, A) \land parentOf(B, C) \land parentOf(C, D) \land male(A)
                 \rightarrow uncleOf(A.D)$
 979
         L11: \sigma(B, C) = 1.11 L11: \sigma(B, A) = 1.11 L11: \sigma(B, C) = 1.11
 980
 981
         L12: $\forall A,B,C,D,E: parentOf(B, A) \land parentOf(B, C) \land parentOf(C, D) \land
 982
 983
                 parentOf(D, E) \land male(A) \rightarrow greatUncleOf(A,E)$
         L13: $\forall A,B,C,D,E,F: parentOf(B, A) \land parentOf(C, B) \land parentOf(C, D) \land parentOf(D, E) \land parentOf(E, F) \land female(A) \rightarrow secondAuntOf(A,F)$
 984
 985
         L14: $\forall A,B,C,D,E,F: parentOf(B, A) \land parentOf(C, B) \land parentOf(C, D) \land
 986
                 parentOf(D, E) \land parentOf(E, F) \land male(A) \rightarrow secondUncleOf(A,F)$
 987
         L15: $\forall A,B,C,D,E: parentOf(B, A) \land parentOf(C, B) \land parentOf(C, D) \land
 988
                 parentOf(D, E) \land female(A) \rightarrow girlCousinOf(A,E)$
 989
 990
         L16: $\forall A,B,C,D,E: parentOf(B, A) \land parentOf(C, B) \land parentOf(C, D) \land
                 parentOf(D, E) \land male(A) \rightarrow boyCousinOf(A,E)$
 991
         L17: \frac{A,B,C,D,E,F,G}{D,E,F,G} parentOf(B, A) \land parentOf(C, B) \land parentOf(D, C) \land
 992
 993
                 parentOf(D, E) \land parentOf(E, F) \land parentOf(F, G) \land female(A) \rightarrow
                 girlSecondCousinOf(A,G)$
 994
         L18: \frac{A,B,C,D,E,F,G}{B}: parentOf(B, A) \frac{A}{B} value parentOf(C, B) \frac{A,B,C,D,E,F,G}{B}:
 995
                parentOf(D, E) \land parentOf(E, F) \land parentOf(F, G) \land male(A) \rightarrow
 996
                 boySecondCousinOf(A,G)$
 997
 998
         L19: $\forall A,B,C,D,E,F: parentOf(B, A) \land parentOf(C, B) \land parentOf(D, C) \land
                 parentOf(D, E) \land parentOf(E, F) \land female(A) \rightarrow
 999
                 girlFirstCousinOnceRemovedOf(A,F)$
1000
1001
         L20: $\forall A,B,C,D,E,F: parentOf(B, A) \land parentOf(C, B) \land parentOf(D, C) \land
                 parentOf(D, E) \land parentOf(E, F) \land male(A) \rightarrow boyFirstCousinOnceRemovedOf
1002
         L21: $\forall A,B: parentOf(B, A) \land female(A) \rightarrow daughterOf(A,B)$
1004
         L22: $\forall A,B: parentOf(B, A) \land male(A) \rightarrow sonOf(A,B)$
1005
1006
         L23: $\forall A,B,C: parentOf(B, A) \land parentOf(C, B) \land female(A) \rightarrow
                 granddaughterOf(A,C)$
1007
         L24: $\forall A,B,C: parentOf(B, A) \land parentOf(C, B) \land male(A) \rightarrow grandsonOf(
1008
                A.C)$
1009
1010
         L25: $\forall A,B,C,D: parentOf(B, A) \land parentOf(C, B) \land parentOf(D, C) \land female(A
1011
                 ) \rightarrow greatGranddaughterOf(A,D)$
         L26: $\forall A,B,C,D: parentOf(B, A) \land parentOf(C, B) \land parentOf(D, C) \land male(A)
1012
                 \rightarrow greatGrandsonOf(A,D)$
1013
         L27: $\forall A,B,C,D: parentOf(B, A) \land parentOf(C, B) \land parentOf(C, D) \land female(A
1014
1015
                 ) \rightarrow nieceOf(A,D)$
         L28: $\forall A,B,C,D: parentOf(B, A) \land parentOf(C, B) \land parentOf(C, D) \land male(A)
1016
1017
                 \rightarrow nephewOf(A,D)$
1018
         Facts:
1019
         F1: female(Laura)
1020
1021
         F2: male(Elias)
         F3: male(Fabian)
1022
        F4: female(Claudia)
1023
1024
         F5: female(Elena)
1025
        F6: male(Thomas)
        F7: female(Amelie)
         F8: female(Luisa)
1027
1028
        F9: male(Patrick)
```

```
F10: female(Emilia)
1029
     F11: male(Samuel)
1030
     F12: female(Alina)
     F13: male(Jonathan)
1032
1033
     F14: male(Philipp)
1034
     F15: male(Nico)
     F16: male(David)
1035
      F17: female(Emily)
1036
     F18: male(Konstantin)
1037
1038
     F19: male(Florian)
1039
     F20: female(Helga)
1040
     F21: female(Nina)
     F22: female(Lea)
1041
1042
     F23: male(Felix)
1043
     F24: female(Leonie)
1044
     F25: male(Stefan)
1045
     F26: male(Gabriel)
      F27: male(Tobias)
1046
     F28: parentOf(Laura, Fabian)
1047
1048
     F29: parentOf(Laura, Felix)
1049
      F30: parentOf(Laura, Claudia)
     F31: parentOf(Elias, Fabian)
1050
     F32: parentOf(Elias, Felix)
     F33: parentOf(Elias, Claudia)
1052
1053
     F34: parentOf(Alina, David)
     F35: parentOf(Alina, Lea)
     F36: parentOf(Nico, David)
1055
1056
     F37: parentOf(Nico, Lea)
     F38: parentOf(Emily, Nico)
1057
     F39: parentOf(Konstantin, Nico)
1058
      F40: parentOf(Fabian, Thomas)
1059
     F41: parentOf(Fabian, Amelie)
1060
     F42: parentOf(Nina, Tobias)
     F43: parentOf(Leonie, Emily)
1062
1063
     F44: parentOf(Stefan, Emily)
      F45: parentOf(Gabriel, Tobias)
     F46: parentOf(Elena, Thomas)
1065
1066
      F47: parentOf(Elena, Amelie)
     F48: parentOf(Thomas, Helga)
1067
     F49: parentOf(Thomas, Nina)
1068
1069
      F50: parentOf(Thomas, Patrick)
     F51: parentOf(Luisa, Helga)
1070
     F52: parentOf(Luisa, Nina)
     F53: parentOf(Luisa, Patrick)
1072
1073
     F54: parentOf(Patrick, Samuel)
1074
     F55: parentOf(Patrick, Alina)
     F56: parentOf(Patrick, Jonathan)
1075
1076
     F57: parentOf(Patrick, Philipp)
1077
      F58: parentOf(Patrick, Florian)
     F59: parentOf(Emilia, Samuel)
1078
1079
      F60: parentOf(Emilia, Alina)
      F61: parentOf(Emilia, Jonathan)
1080
1081
      F62: parentOf(Emilia, Philipp)
1082
      F63: parentOf(Emilia, Florian)
1083
1084
      Unknown fact: boyCousinOf(Tobias, David)
```

B.1.2 Natural language representations

```
L1: If B is parent of A and B is parent of C and A is female, then A is sister of D.
1087
1088
     L2: If B is parent of A and B is parent of C and A is male, then A is brother of D.
1089
     L3: If A is parent of B and A is female, then A is mother of C.
     L4: If A is parent of B and A is male, then A is father of C.
1090
     L5: If A is parent of B and B is parent of C and A is female, then A is grandmother of D.
1091
1092
     L6: If A is parent of B and B is parent of C and A is male, then A is grandfather of D.
1093
     L7: If A is parent of B and B is parent of C and C is parent of D and A is female, then A is
          greatGrandmother of E.
     L8: \overline{\text{If}} A is parent of B and B is parent of C and C is parent of D and A is male, then A is
1095
1096
         greatGrandfather of E.
```

```
L9: If B is parent of A and B is parent of C and C is parent of D and A is female, then A is
1097
1098
          aunt of E.
     L10: If B is parent of A and B is parent of C and C is parent of D and A is male, then A is
1099
1100
          uncle of E.
1101
     L11: If B is parent of A and B is parent of C and C is parent of D and D is parent of E and A
1102
          is female, then A is greatAunt of F.
     L12: If B is parent of A and B is parent of C and C is parent of D and D is parent of E and A
1103
           is male, then A is greatUncle of F.
1104
     L13: If B is parent of A and C is parent of B and C is parent of D and D is parent of E and E
1105
1106
           is parent of F and A is female, then A is secondAunt of G.
1107
     L14: If B is parent of A and C is parent of B and C is parent of D and D is parent of E and E
          is parent of F and A is male, then A is secondUncle of G.
1108
     L15: If B is parent of A and C is parent of B and C is parent of D and D is parent of E and A
1109
1110
           is female, then A is girlCousin of F.
1111
     L16: If B is parent of A and C is parent of B and C is parent of D and D is parent of E and A
1112
          is male, then A is boyCousin of F.
1113
     L17: If B is parent of A and C is parent of B and D is parent of C and D is parent of E and E
           is parent of F and F is parent of G and A is female, then A is girlSecondCousin of H.
1114
     L18: If B is parent of A and C is parent of B and D is parent of C and D is parent of E and E
1115
1116
          is parent of F and F is parent of G and A is male, then A is boySecondCousin of H.
1117
     L19: If B is parent of A and C is parent of B and D is parent of C and D is parent of E and E
           is parent of F and A is female, then A is girlFirstCousinOnceRemoved of G.
1118
     L20: If B is parent of A and C is parent of B and D is parent of C and D is parent of E and E
          is parent of F and A is male, then A is boyFirstCousinOnceRemoved of {\tt G.}
1120
1121
     L21: If B is parent of A and A is female, then A is daughter of C.
     L22: If B is parent of A and A is male, then A is son of C.
     L23: If B is parent of A and C is parent of B and A is female, then A is granddaughter of D.
1123
     L24: If B is parent of A and C is parent of B and A is male, then A is grandson of D.
1124
     L25: If B is parent of A and C is parent of B and D is parent of C and A is female, then A is
1125
          greatGranddaughter of E.
1126
     L26: If B is parent of A and C is parent of B and D is parent of C and A is male, then A is
1127
          greatGrandson of E.
1128
1129
     L27: If B is parent of A and C is parent of B and C is parent of D and A is female, then A is
          niece of E.
1130
     L28: If B is parent of A and C is parent of B and C is parent of D and A is male, then A is
1131
1132
          nephew of E.
1133
1134
     Facts:
     F1: Laura is female.
1135
1136
     F2: Elias is male.
1137
     F3: Fabian is male
     F4: Claudia is female.
1138
     F5: Elena is female.
     F6: Thomas is male.
1140
1141
     F7: Amelie is female.
1142
     F8: Luisa is female.
     F9: Patrick is male.
1143
1144
     F10: Emilia is female.
1145 F11: Samuel is male.
1146 F12: Alina is female.
1147
     F13: Jonathan is male.
     F14: Philipp is male.
1148
1149
     F15: Nico is male.
1150
     F16: David is male.
     F17: Emily is female.
1151
     F18: Konstantin is male.
     F19: Florian is male.
1153
1154
     F20: Helga is female.
     F21: Nina is female.
1155
     F22: Lea is female.
1156
1157
     F23: Felix is male.
     F24: Leonie is female.
1158
1159
     F25: Stefan is male.
     F26: Gabriel is male.
1160
     F27: Tobias is male.
1161
     F28: Laura is parent of Fabian.
     F29: Laura is parent of Felix.
1163
     F30: Laura is parent of Claudia.
1164
     F31: Elias is parent of Fabian.
     F32: Elias is parent of Felix.
1166
1167
     F33: Elias is parent of Claudia.
1168 F34: Alina is parent of David.
```

```
F35: Alina is parent of Lea.
1169
1170
     F36: Nico is parent of David.
     F37: Nico is parent of Lea.
     F38: Emily is parent of Nico.
1172
     F39: Konstantin is parent of Nico.
1173
1174
     F40: Fabian is parent of Thomas.
     F41: Fabian is parent of Amelie.
1175
     F42: Nina is parent of Tobias.
1176
     F43: Leonie is parent of Emily.
1177
     F44: Stefan is parent of Emily.
1178
     F45: Gabriel is parent of Tobias.
1179
1180
     F46: Elena is parent of Thomas.
     F47: Elena is parent of Amelie.
     F48: Thomas is parent of Helga.
1182
     F49: Thomas is parent of Nina.
1183
     F50: Thomas is parent of Patrick.
1184
1185
     F51: Luisa is parent of Helga.
1186
     F52: Luisa is parent of Nina.
     F53: Luisa is parent of Patrick.
1187
     F54: Patrick is parent of Samuel.
1188
     F55: Patrick is parent of Alina.
1189
1190
     F56: Patrick is parent of Jonathan.
1191
     F57: Patrick is parent of Philipp.
     F58: Patrick is parent of Florian.
1192
1193
     F59: Emilia is parent of Samuel.
     F60: Emilia is parent of Alina.
1194
     F61: Emilia is parent of Jonathan.
1195
     F62: Emilia is parent of Philipp.
1196
     F63: Emilia is parent of Florian.
1197
1198
1199
     Unknown fact: Gabriel is uncle of Lea.
```

1200 B.2 Symbolization

1201 B.2.1 Logic language representations

```
Logical rules:
1202
          L1: $\forall A,B,C: r3(B, A) \land r3(B, C) \land r2(A) \rightarrow r4(A, C)$
1203
          L2: $\forall A,B,C: r3(B, A) \land r3(B, C) \land r1(A) \rightarrow r5(A, C)$
1204
1205
          L3: \frac{A}{B}: r3(A, B) \frac{A}{B} r2(A) \frac{A}{B}
1206
          L4: $\forall A,B: r3(A, B) \land r1(A) \rightarrow r7(A, B)$
          L5: $\forall A,B,C: r3(A, B) \land r3(B, C) \land r2(A) \rightarrow r8(A, C)$
1207
          L6: $\forall A,B,C: r3(A, B) \land r3(B, C) \land r1(A) \rightarrow r9(A, C)$
1209
          L7: $\forall A,B,C,D: r3(A, B) \land r3(B, C) \land r3(C, D) \land r2(A) \rightarrow r10(A, D)
1210
          L8: $\forall A,B,C,D: r3(A, B) \land r3(B, C) \land r3(C, D) \land r1(A) \rightarrow r11(A, D)
1211
                    $
1212
           L9: $\forall A,B,C,D: r3(B, A) \land r3(B, C) \land r3(C, D) \land r2(A) \rightarrow r12(A, D)
1213
1214
          L10: $\forall A,B,C,D: r3(B, A) \land r3(B, C) \land r3(C, D) \land r1(A) \rightarrow r13(A, D
1215
1216
          L11: $\forall A,B,C,D,E: r3(B, A) \land r3(B, C) \land r3(C, D) \land r3(D, E) \land r2(A) \
1217
1218
                    rightarrow r14(A, E)$
1219
           L12: $\forall A,B,C,D,E: r3(B, A) \land r3(B, C) \land r3(C, D) \land r3(D, E) \land r1(A) \
1220
                    rightarrow r15(A, E)$
           L13: $\forall A,B,C,D,E,F: r3(B, A) \land r3(C, B) \land r3(C, D) \land r3(D, E) \land r3(E, F
                   ) \land r2(A) \rightarrow r16(A, F)$
1222
1223
          L14: $\forall A,B,C,D,E,F: r3(B, A) \land r3(C, B) \land r3(C, D) \land r3(D, E) \land r3(E, F
                   ) \land r1(A) \rightarrow r17(A, F)$
1224
          L15: \frac{1}{5} A,B,C,D,E: r3(B, A) \land r3(C, B) \land r3(C, D) \land r3(D, E) \land r2(A) \
1225
1226
                    rightarrow r18(A, E)$
          L16: forall A,B,C,D,E: r3(B, A) \ r3(C, B) \ r3(C, D) \ r3(D, E) \ r1(A) \ r3(B, E) \ r3(B, E) \ r3(B, E) \ r3(B, E) \ r3(E, E) \ 
1227
                    rightarrow r19(A, E)$
1228
          L17: $\forall A,B,C,D,E,F,G: r3(B, A) \land r3(C, B) \land r3(D, C) \land r3(D, E) \land r3(E,
1229
                     F) \land r3(F, G) \land r2(A) \rightarrow r20(A, G)$
1230
          L18: $\forall A,B,C,D,E,F,G: r3(B, A) \land r3(C, B) \land r3(D, C) \land r3(D, E) \land r3(E,
1231
1232
                     F) \land r3(F, G) \land r1(A) \rightarrow r21(A, G)$
          L19: $\forall A,B,C,D,E,F: r3(B, A) \land r3(C, B) \land r3(D, C) \land r3(D, E) \land r3(E, F
1233
1234
                   ) \land r2(A) \rightarrow r22(A, F)$
          L20: $\forall A,B,C,D,E,F: r3(B, A) \land r3(C, B) \land r3(D, C) \land r3(D, E) \land r3(E, F
1235
           ) \land r1(A) \rightarrow r23(A, F)$
1236
```

```
L21: forall A,B: r3(B, A) \arrow r2(A) \rightarrow r24(A, B) 
L22: <math>forall A,B: r3(B, A) \arrow r1(A) \rightarrow r25(A, B)
1237
1238
      L23: $\forall A,B,C: r3(B, A) \land r3(C, B) \land r2(A) \rightarrow r26(A, C)$
      L24: $\forall A,B,C: r3(B, A) \land r3(C, B) \land r1(A) \rightarrow r27(A, C)$
1240
1241
      L25: $\forall A,B,C,D: r3(B, A) \land r3(C, B) \land r3(D, C) \land r2(A) \rightarrow r28(A, D
1242
      L26: $\forall A,B,C,D: r3(B, A) \land r3(C, B) \land r3(D, C) \land r1(A) \rightarrow r29(A, D
1243
1244
           )$
      L27: $\forall A,B,C,D: r3(B, A) \land r3(C, B) \land r3(C, D) \land r2(A) \rightarrow r30(A, D
1245
1246
           )$
1247
      L28: $\forall A,B,C,D: r3(B, A) \land r3(C, B) \land r3(C, D) \land r1(A) \rightarrow r31(A, D
1248
           )$
1249
1250
      Facts:
1251
      F1: $r2$(Laura)
      F2: $r1$(Elias)
1252
1253
      F3: $r1$(Fabian)
      F4: $r2$(Claudia)
1254
      F5: $r2$(Elena)
1255
1256
      F6: $r1$(Thomas)
1257
      F7: $r2$(Amelie)
     F8: $r2$(Luisa)
1258
      F9: $r1$(Patrick)
      F10: $r2$(Emilia)
1260
1261
     F11: $r1$(Samuel)
     F12: $r2$(Alina)
     F13: $r1$(Jonathan)
1263
1264
      F14: $r1$(Philipp)
     F15: $r1$(Nico)
1265
1266
     F16: $r1$(David)
      F17: $r2$(Emily)
1267
     F18: $r1$(Konstantin)
1268
1269
      F19: $r1$(Florian)
      F20: $r2$(Helga)
1270
1271
     F21: $r2$(Nina)
      F22: $r2$(Lea)
1272
      F23: $r1$(Felix)
1273
1274
      F24: $r2$(Leonie)
     F25: $r1$(Stefan)
1275
1276
     F26: $r1$(Gabriel)
1277
      F27: $r1$(Tobias)
     F28: $r3$(Laura, Fabian)
1278
      F29: $r3$(Laura, Felix)
      F30: $r3$(Laura, Claudia)
1280
1281
      F31: $r3$(Elias, Fabian)
      F32: $r3$(Elias, Felix)
1282
      F33: $r3$(Elias, Claudia)
1283
     F34: $r3$(Alina, David)
1284
1285
      F35: $r3$(Alina, Lea)
     F36: $r3$(Nico, David)
1286
1287
      F37: $r3$(Nico, Lea)
     F38: $r3$(Emily, Nico)
1288
1289
     F39: $r3$(Konstantin, Nico)
1290
      F40: $r3$(Fabian, Thomas)
      F41: $r3$(Fabian, Amelie)
1291
      F42: $r3$(Nina, Tobias)
      F43: $r3$(Leonie, Emily)
1293
1294
      F44: $r3$(Stefan, Emily)
      F45: $r3$(Gabriel, Tobias)
1295
      F46: $r3$(Elena, Thomas)
1296
1297
      F47: $r3$(Elena, Amelie)
      F48: $r3$(Thomas, Helga)
1298
1299
      F49: $r3$(Thomas, Nina)
      F50: $r3$(Thomas, Patrick)
1300
      F51: $r3$(Luisa, Helga)
1301
1302
      F52: $r3$(Luisa, Nina)
      F53: $r3$(Luisa, Patrick)
1303
1304
      F54: $r3$(Patrick, Samuel)
1305
      F55: $r3$(Patrick, Alina)
     F56: $r3$(Patrick, Jonathan)
F57: $r3$(Patrick, Philipp)
1306
1307
1308 F58: $r3$(Patrick, Florian)
```

```
1309 F59: $r3$(Emilia, Samuel)
1310 F60: $r3$(Emilia, Alina)
1311 F61: $r3$(Emilia, Jonathan)
1312 F62: $r3$(Emilia, Philipp)
1313 F63: $r3$(Emilia, Florian)
1314
1315 Unknown fact: $r9$(Thomas, Claudia)
```

1316 B.2.2 Natural language representations:

```
1317
     Logical rules:
     L1: If B is $r3$ of A and B is $r3$ of C and A is $r2$, then A is $r4$ of D.
1318
     L2: If B is $r3$ of A and B is $r3$ of C and A is $r1$, then A is $r5$ of D.
     L3: If A is $r3$ of B and A is $r2$, then A is $r6$ of C.
1320
1321
     L4: If A is r3 of B and A is r1, then A is r7 of C.
     L5: If A is $r3$ of B and B is $r3$ of C and A is $r2$, then A is $r8$ of D.
     L6: If A is r3 of B and B is r3 of C and A is r1, then A is r9 of D.
1323
1324
     L7: If A is $r3$ of B and B is $r3$ of C and C is $r3$ of D and A is $r2$, then A is $r10$ of
1325
          E
1326
     L8: If A is $r3$ of B and B is $r3$ of C and C is $r3$ of D and A is $r1$, then A is $r11$ of
1327
          Ε.
     L9: If B is $r3$ of A and B is $r3$ of C and C is $r3$ of D and A is $r2$, then A is $r12$ of
1328
1329
          E.,
1330
     L10: If B is $r3$ of A and B is $r3$ of C and C is $r3$ of D and A is $r1$, then A is $r13$ of
1331
           E.,
1332
      L11: If B is $r3$ of A and B is $r3$ of C and C is $r3$ of D and D is $r3$ of E and A is $r2$,
           then A is $r14$ of F.
1333
1334
     L12: If B is $r3$ of A and B is $r3$ of C and C is $r3$ of D and D is $r3$ of E and A is $r1$,
           then A is $r15$ of F.
1335
1336
     L13: If B is $r3$ of A and C is $r3$ of B and C is $r3$ of D and D is $r3$ of E and E is $r3$
1337
          of F and A is $r2$, then A is $r16$ of G.
     L14: If B is $r3$ of A and C is $r3$ of B and C is $r3$ of D and D is $r3$ of E and E is $r3$
1338
1339
          of F and A is $r1$, then A is $r17$ of G.
     L15: If B is $r3$ of A and C is $r3$ of B and C is $r3$ of D and D is $r3$ of E and A is $r2$,
1340
1341
           then A is $r18$ of F.
1342
     L16: If B is $r3$ of A and C is $r3$ of B and C is $r3$ of D and D is $r3$ of E and A is $r1$,
           then A is $r19$ of F.
1343
     L17: If B is $r3$ of A and C is $r3$ of B and D is $r3$ of C and D is $r3$ of E and E is $r3$
1344
          of F and F is $r3$ of G and A is $r2$, then A is $r20$ of H.
1345
1346
     L18: If B is $r3$ of A and C is $r3$ of B and D is $r3$ of C and D is $r3$ of E and E is $r3$
1347
          of F and F is $r3$ of G and A is $r1$, then A is $r21$ of H.
     L19: If B is $r3$ of A and C is $r3$ of B and D is $r3$ of C and D is $r3$ of E and E is $r3$
1348
          of F and A is $r2$, then A is $r22$ of G.
1349
1350
     L20: If B is $r3$ of A and C is $r3$ of B and D is $r3$ of C and D is $r3$ of E and E is $r3$
          of F and A is $r1$, then A is $r23$ of G.
1351
     L21: If B is $r3$ of A and A is $r2$, then A is $r24$ of C.
1352
     L22: If B is $r3$ of A and A is $r1$, then A is $r25$ of C.
1353
1354
     L23: If B is $r3$ of A and C is $r3$ of B and A is $r2$, then A is $r26$ of D.
1355
     L24: If B is $r3$ of A and C is $r3$ of B and A is $r1$, then A is $r27$ of D.
     L25: If B is $r3$ of A and C is $r3$ of B and D is $r3$ of C and A is $r2$, then A is $r28$ of
1356
1357
           Ε.
     L26: If B is $r3$ of A and C is $r3$ of B and D is $r3$ of C and A is $r1$, then A is $r29$ of
1358
1359
           Ε.
1360
      L27: If B is $r3$ of A and C is $r3$ of B and C is $r3$ of D and A is $r2$, then A is $r30$ of
1361
           Ε.
      L28: If B is $r3$ of A and C is $r3$ of B and C is $r3$ of D and A is $r1$, then A is $r31$ of
1362
           E.
1363
1364
1365
     Facts:
1366
     F1: Laura is $r2$.
1367
     F2: Elias is $r1$.
     F3: Fabian is $r1$.
1368
     F4: Claudia is $r2$.
1369
1370
     F5: Elena is $r2$.
1371
     F6: Thomas is $r1$.
     F7: Amelie is $r2$.
1372
1373
     F8: Luisa is $r2$.
1374
     F9: Patrick is $r1$.
1375 F10: Emilia is $r2$.
     F11: Samuel is $r1$.
1376
     F12: Alina is $r2$.
1377
```

```
1378
     F13: Jonathan is $r1$.
1379
     F14: Philipp is $r1$.
1380 F15: Nico is $r1$.
     F16: David is $r1$.
1381
1382
     F17: Emily is $r2$.
1383 F18: Konstantin is $r1$.
     F19: Florian is $r1$.
1384
1385
     F20: Helga is $r2$.
     F21: Nina is $r2$.
1386
1387
     F22: Lea is $r2$.
1388
     F23: Felix is $r1$.
1389
     F24: Leonie is $r2$.
     F25: Stefan is $r1$.
1390
1391
     F26: Gabriel is $r1$.
1392
     F27: Tobias is $r1$.
     F28: Laura is $r3$ of Fabian.
1393
1394
     F29: Laura is $r3$ of Felix.
     F30: Laura is $r3$ of Claudia.
1395
     F31: Elias is $r3$ of Fabian.
1396
1397
     F32: Elias is $r3$ of Felix.
1398
     F33: Elias is $r3$ of Claudia.
     F34: Alina is $r3$ of David.
1399
     F35: Alina is $r3$ of Lea.
     F36: Nico is $r3$ of David.
1401
1402
     F37: Nico is $r3$ of Lea.
1403 F38: Emily is $r3$ of Nico.
     F39: Konstantin is $r3$ of Nico.
1404
1405
     F40: Fabian is $r3$ of Thomas.
     F41: Fabian is $r3$ of Amelie.
1406
1407
     F42: Nina is $r3$ of Tobias.
     F43: Leonie is $r3$ of Emily.
1408
     F44: Stefan is $r3$ of Emily.
1409
1410
     F45: Gabriel is $r3$ of Tobias.
     F46: Elena is $r3$ of Thomas.
1411
1412
     F47: Elena is $r3$ of Amelie.
     F48: Thomas is $r3$ of Helga.
1413
     F49: Thomas is $r3$ of Nina.
1414
1415
     F50: Thomas is $r3$ of Patrick.
     F51: Luisa is $r3$ of Helga.
1416
1417
     F52: Luisa is $r3$ of Nina.
1418
     F53: Luisa is $r3$ of Patrick.
     F54: Patrick is $r3$ of Samuel.
1419
     F55: Patrick is $r3$ of Alina.
     F56: Patrick is $r3$ of Jonathan.
1421
1422
     F57: Patrick is $r3$ of Philipp.
1423
     F58: Patrick is $r3$ of Florian.
     F59: Emilia is $r3$ of Samuel.
1424
1425
     F60: Emilia is $r3$ of Alina.
1426
     F61: Emilia is $r3$ of Jonathan.
     F62: Emilia is $r3$ of Philipp.
1427
1428
     F63: Emilia is $r3$ of Florian.
1429
     Unknown fact: Nico is $r27$ of Stefan.
1430
```

B.3 Semantics of removing rule setting

1431

```
I will provide a set of facts. Please predict True/False of the unknown fact based on given
1432
1433
          facts.
1434
     Facts:
1435
     F1: Laura is female.
     F2: Elias is male.
1436
     F3: Fabian is male.
1437
1438 F4: Claudia is female.
1439
     F5: Elena is female.
1440
     F6: Thomas is male.
1441
     F7: Amelie is female.
1442
     F8: Luisa is female.
1443 F9: Patrick is male.
1444 F10: Emilia is female.
     F11: Samuel is male.
1445
1446
     F12: Alina is female.
```

```
F13: Jonathan is male.
1447
1448
     F14: Philipp is male.
1449 F15: Nico is male.
     F16: David is male.
1450
1451
     F17: Emily is female.
1452
    F18: Konstantin is male.
     F19: Florian is male.
1453
     F20: Helga is female.
1454
     F21: Nina is female.
1455
1456
     F22: Lea is female.
1457
     F23: Felix is male.
1458
     F24: Leonie is female.
     F25: Stefan is male.
1459
1460
     F26: Gabriel is male.
1461
     F27: Tobias is male.
     F28: Laura is parent of Fabian.
1462
1463
     F29: Laura is parent of Felix.
     F30: Laura is parent of Claudia.
     F31: Elias is parent of Fabian.
1465
1466
     F32: Elias is parent of Felix.
1467
     F33: Elias is parent of Claudia.
1468
     F34: Alina is parent of David.
     F35: Alina is parent of Lea.
     F36: Nico is parent of David.
1470
1471
     F37: Nico is parent of Lea.
     F38: Emily is parent of Nico.
1472
     F39: Konstantin is parent of Nico.
1473
1474
     F40: Fabian is parent of Thomas.
     F41: Fabian is parent of Amelie.
1475
1476
     F42: Nina is parent of Tobias.
1477
     F43: Leonie is parent of Emily.
     F44: Stefan is parent of Emily.
1478
1479
     F45: Gabriel is parent of Tobias.
     F46: Elena is parent of Thomas.
1480
1481
     F47: Elena is parent of Amelie.
     F48: Thomas is parent of Helga.
     F49: Thomas is parent of Nina.
1483
1484
     F50: Thomas is parent of Patrick.
     F51: Luisa is parent of Helga.
1485
1486
     F52: Luisa is parent of Nina.
1487
     F53: Luisa is parent of Patrick.
     F54: Patrick is parent of Samuel.
1488
     F55: Patrick is parent of Alina.
1490
     F56: Patrick is parent of Jonathan.
1491
     F57: Patrick is parent of Philipp.
1492
     F58: Patrick is parent of Florian.
     F59: Emilia is parent of Samuel.
1493
1494
     F60: Emilia is parent of Alina.
     F61: Emilia is parent of Jonathan.
1495
     F62: Emilia is parent of Philipp.
1496
1497
     F63: Emilia is parent of Florian.
1498
1499
     Unknown fact: Jonathan is aunt of Thomas.
1500
     The answer (True or False) is:
```

C Examples of ProofWriter

In this section, we provide examples of deduction experiments conducted on the ProofWriter Depth-1 dataset. We present examples for both the *Semantics* and *Symbols* settings.

C.1 Semantics

1501

1504

```
The bear likes the dog.

The cow is round.

The cow likes the bear.

The cow needs the bear.

The dog needs the squirrel.

The dog sees the cow.

The squirrel needs the dog.
```

```
1512 If someone is round then they like the squirrel.
1513 If the bear is round and the bear likes the squirrel then the squirrel needs the bear.
1514 If the cow needs the dog then the cow is cold.
1515 Does it imply that the statement "The cow likes the squirrel." is True?
1516 The bear likes the dog.
     The cow is round.
1517
     The cow likes the bear.
1518
1519 The cow needs the bear.
     The dog needs the squirrel.
     The dog sees the cow.
1521
1522 The squirrel needs the dog.
1523
     If someone is round then they like the squirrel.
1524
     If the bear is round and the bear likes the squirrel then the squirrel needs the bear.
1525 If the cow needs the dog then the cow is cold.
1526 Does it imply that the statement "The cow does not like the squirrel." is True?
1527 Bob is blue.
     Erin is quiet.
1528
1529 Fiona is cold.
1530 Harry is cold.
1531 All quiet things are blue.
1532 If Harry is blue then Harry is not young.
1533 Blue things are young.
     Blue, round things are cold.
1534
1535
     If something is blue and not red then it is round.
     If something is young then it is white.
1536
1537
     If Erin is red and Erin is not round then Erin is young.
     If Erin is red and Erin is not cold then Erin is white.
     Does it imply that the statement "Erin is white" is True?
1539
    Answer with only True or False. The answer is:
1540
1541 The bear likes the dog.
     The cow is round.
1542
1543
     The cow likes the bear.
1544 The cow needs the bear.
1545 The dog needs the squirrel.
     The dog sees the cow.
1546
1547
     The squirrel needs the dog.
1548 If someone is round then they like the squirrel.
     If the bear is round and the bear likes the squirrel then the squirrel needs the bear.
1550
     If the cow needs the dog then the cow is cold.
     Does it imply that the statement "The cow likes the squirrel." is True?
    C.2 Symbols
1552
     The e4 likes the e5.
1554
     The e14 is e2.
1555 The e14 likes the e4.
     The e14 needs the e4.
1556
1557 The e5 needs the e26.
1558 The e5 sees the e14.
1559
     The e26 needs the e5.
1560
     If someone is e2 then they like the e26.
    If the e4 is e2 and the e4 likes the e26 then the e26 needs the e4.
1562
     If the e14 needs the e5 then the e14 is e1.
1563
     Does it imply that the statement "The e14 likes the e26." is True?
     The e27 is e7.
     The e27 is e15.
1565
1566
     The e30 does not chase the e27.
1567
     The e30 eats the e27.
     The e30 is e1.
1568
     The e30 is e15.
1569
1570
     The e30 visits the e27.
1571
    If something visits the e27 then the e27 does not visit the e30.
1572 If something is e1 and e15 then it visits the e30.
1573 Does it imply that the statement "The e30 visits the e30." is True?
```

```
The e27 is e7.
1574
     The e27 is e15.
1575
     The e30 does not chase the e27.
     The e30 eats the e27.
1577
1578
     The e30 is e1.
1579
     The e30 is e15.
1580
     The e30 visits the e27.
      If something visits the e27 then the e27 does not visit the e30.
1581
     If something is e1 and e15 then it visits the e30.
1582
     Does it imply that the statement "The e30 visits the e30." is True?
1583
```

D Different Zero-Shot prompting

1585 We try different Zero-Shot prompts:

Here are some facts: {basic facts}

The answer (YES or NO) is:

Does it imply that the statement "{unknown fact}" is True?

1584

1605

1606

1607

1610

1611

```
(1)
1586
     I will provide a set of logical rules L1 to L{number of rules} and facts F1 to F{number of
1587
          basic facts}. Please select one single logical rule from L1 to L{number of rules} and a
1588
           few facts from F1 to F{number of basic facts} to predict True/False of the unknown fact
1589
1590
          using deductive reasoning.
1591
     Logical rules: {rules}
1592
     Facts: {basic facts}
     Unknown fact: {unknown fact}
1593
     The answer (True or False) is:
1594
     (2)
1595
1596
     I will provide a set of logical rules L1 to L{number of rules} and facts F1 to F{number of
           basic facts}. Please predict True/False of the unknown fact using deductive reasoning.
1597
1598
     Logical rules: {rules}
1599
     Facts: {basic facts}
1600
     Unknown fact: {unknown fact}
     The answer (True or False) is:
1601
     (3)
1602
     Given a set of rules and facts, you have to reason whether a statement is True or False.
1603
     Here are some rules: {rules}
1604
```

The results of the three prompts in the Zero-Shot setting are presented in Table 5. Among the three prompts, we select the one that achieves the best performance as our Zero-Shot prompt.

Table 5: Different Zero-Shot Prompts of deductive reasoning. Results are in %.

1	_	_
nromntl	nromnt2	nromnt3

	prompt1	prompt2	prompt3
KG_1	54.5	51.5	53.8

E Comparison of memorization abilities of neural-based and symbolic-based methods

We compare fine-tuned language models with the deterministic graph DB Neo4J to explore the 1612 memorization abilities of neural-based and symbolic-based methods. Language models can implicitly store and retrieve facts as "knowledge bases" within their neural parameters. They are trained on a 1614 snapshot of data and may not have access to the latest or most accurate information. In order to update 1615 1616 or add facts, specific model parameters need to be modified, or the model needs to be fine-tuned with new data. In contrast, symbolic knowledge graphs can directly add or update individual triplets, 1617 making it easier to incorporate new information. Our comparison affirms the huge advantage of 1618 using KGs/external DBs to update knowledge rather than finetuning, aligning with the recent trend of 1619 retrieval-based LLM. 1620

F Introduction of Neo4j

The Symbolic Tree is also a knowledge graph dataset. We conduct a comparison between the memorization abilities of a popular graph database, **Neo4j**, and LLMs, **LLaMA-7B**. Neo4j is a widely used graph database system that provides convenient operations such as querying, inserting, deleting, and revising knowledge graphs. For our comparison, we deployed Neo4j on a high-performance server equipped with 2 Intel(R) Xeon(R) Platinum 8380 CPUs, each with 40 cores and 80 threads. The server has 512GB of memory and 4x1.8T NVME SSD disks.

To ensure a fair comparison, we configured Neo4j with a pre-stored knowledge base that has a comparable disk space size to the LLaMA language model. Specifically, we used the Freebase dataset for Neo4j, which occupies approximately 30GB of disk space after preprocessing. For the language model, we used LLaMA-7B, which requires about 14GB of disk space. By comparing the performance of Neo4j and LLaMA-7B in terms of their memorization abilities, we can gain insights into the advantages and limitations of graph databases and language models for storing and retrieving knowledge.

1635 G Task definitions

We define a few tasks to evaluate LLMs' abilities of three kinds of reasoning and memorization.

- deductive reasoning: we use hypothesis classification, i.e., predict the correctness of the hypothesis given the theory where theory consists of basic facts and logical rules, correctness can be true or false, and hypothesis is a predicted fact, which is one of the inferred facts or negative samples. The accuracy is the proportion of correct predictions.
- inductive reasoning: we perform the rule generation task. Given multiple facts with similar patterns and a rule template, the goal is to induce a rule that entails these facts. Specifically, for each relation r, we use basic facts and those inferred facts that contain only relation r as provided facts. The induced rule is generated after filling in the rule template. We test the generated rules against the ground truth rules. If the generated rule matches the ground truth rule exactly, we predict the rule to be correct; otherwise, we predict the rule to be incorrect. The precision is the proportion of correct predictions. Note that considering logical rules maybe not all chain rules (e.g., $r_1(y,x) \wedge r_2(y,z) \rightarrow r_3(x,z)$), we add inverse relation for each relation in order to transform them into chain rules and simplify the rule template (e.g., $r_1^{-1}(x,y) \wedge r_2(y,z) \rightarrow r_3(x,z)$). Furthermore, we provide a rule template for each relation. Take auntOf as example, its rule template can be $\forall x, y, z : \#(x,y) \wedge \#(y,z) \wedge ++(x) \rightarrow auntOf(x,z)$ or "If x is # of y and y is # of z and x is ++, then x is aunt of z.", where # can be parent or parent, parent or parent.
 - Besides, a single rule can be equivalent to multiple rules. For example, the rule $\forall x,y,z:$ parentOf $(x,y) \land \text{parentOf}(y,z) \land \text{gender}(x,\text{female}) \rightarrow \text{GrandmotherOf}(x,z)$ can be represented as $\forall x,y,z:$ parentOf $(x,y) \land \text{parentOf}(y,z) \rightarrow \text{GrandparentOf}(x,z)$, GrandparentOf $(x,z) \land \text{gender}(x,\text{female}) \rightarrow \text{GrandmotherOf}(x,z)$. We conduct the experiments with both rule representations and find single-longer rules perform better than multiple-short rules. Results are presented in Appendix R. Based on these observations and considering the simplicity of induction evaluation, we rewrite all logical rules by including only the parentOf and gender relations in the rule body. This also ensures that each inferred relation is implied by a single logical rule, referred to as $grounding\ truth\ rule$.
- abductive reasoning: We use explanation generation to evaluate abductive reasoning abilities. Given a theory including basic facts and all logical rules, the task is to select specific facts and a logical rule to explain the observation. The observation is chosen from inferred facts. We use Proof Accuracy (PA) as an evaluation metric, i.e., the fraction of examples where the generated proof matches exactly any of the gold proofs.
- memorization: We use a subset of Symbolic Trees to fine-tune the language model. For the symbolic setting, we use r_1, r_2, r_3 to replace the original relations in the semantic setting. Note that the new dataset does not overlap with the old knowledge base of LLMs, ensuring no disambiguation problem and the influence of pre-existing knowledge. When memorizing, we use time, efficiency and forgetting as metrics: time is the cost time of adding/updating facts, efficiency is the MRR (mean reciprocal rank [71]) of facts added/updated, and forgetting is the MRR of the

facts that should not be updated. When evaluating whether a fact has been successfully added or updated, we query LLM with a question about the tail entity and rank the probability of all tokens between all tail entities. The better LLM remembers a triplet, the higher the MRR gets. Note that, there may be more than one entity for each (head, relation) pair. We only consider the rank one of them.

H Implementation of memorization

1679

1693

1715

We selected 1258 triplets from 4 Symbolic Trees to evaluate the effectiveness of adding knowledge. 1680 Following the prompting of Taori et al. [74], we use the head entity and relation as instructions and provide all candidate tails as input. The model's training objective is to autoregress toward the true 1682 tail entities. The detailed prompting is contained in Appendix A.5. In the updating step, we fine-tune 1683 the model on all 620 triplets from the first two trees whose tail entities are randomly flipped to false 1684 ones. Besides the effectiveness of updating, we evaluated the forgetting ratio using the remaining 1685 638 triplets of the least two trees. These triplets have been remembered in the first step and haven't 1686 been updated in the second. Noting that, within each tree, the relationships between entities are 1687 1688 independent, and the entities are distinct. Therefore, we propose that LLM should retain its memory of the previously remembered triplets when updating based on the first two trees. We utilized 4 A100 1689 80G GPUs with batch size 64 for finetuning. The training process involved 100 epochs, employing a 1690 cosine learning rate schedule with an initial learning rate of 2e-5. We run these experiments three 1691 times and recorded their mean and standard MRR. 1692

I Consistency of knowledge base

In the context of updating a knowledge base, it is important to ensure the consistency of the knowledge 1694 base. When revising a fact, it is necessary to update other related facts accordingly to maintain a 1695 coherent and accurate knowledge base. For example, if we have the facts: Alice is Bob's mother, 1696 Amy is Alice's sister, Momo is Lily's sister, and Amy is Bob's aunt, and we update Bob's mother as 1697 Lily, it follows that Bob's aunt should be updated to Momo instead of Amy to maintain consistency. 1698 However, in our experiments, we did not consider the situation of updating related facts in order to 1699 simplify the evaluation process, although it may lead to inconsistencies within the knowledge base. 1700 In real-world scenarios, ensuring consistency is an essential aspect of maintaining the accuracy and reliability of a knowledge base when updating or revising facts.

1703 J Different Few-Shot-CoT examples

- We use three different kinds of demonstrations generation:
- (1) leverage LLMs with the "Let's think step by step" prompt to generate reasoning chains for demonstrations.
- (2) manually designed each demonstration with a question and a reasoning chain. A reasoning chain is composed of a rationale (a series of intermediate reasoning steps) and an expected answer.
- (3) leverage LLMs with the "Let's think step by step" prompt to generate reasoning chains and filter with those generated chains coming with mistakes for demonstrations.
- The results of these approaches are provided in Table 6. From the results, we observe that the Auto-Randomly approach outperforms both the Auto-Filter and Manually-Designed approaches. This suggests that randomly generating reasoning chains is more suitable for LLMs and can generate a diverse range of solutions, resulting in better reasoning performance.

J.1 Few-Shot-CoT examples with Auto-Filter

- 1716 Statement: r8(Elena, Nina)
- Answer: We can use logical rule L5: $\forall A, B, C : r3(A, B) \land r3(B, C) \land r2(A) \rightarrow r8(A, C)$ to
- deduce whether the statement r8 (Elena, Nina) is true or false. From the facts, we know that r3 (Elena,
- Thomas) and r3 (Thomas, Nina). Therefore, we can use L5 with A = Elena, B = Thomas, and
- C = N in a to see if r8(Elena, Nina) is true or false. L5 states that if r3(A, B) and r3(B, C) and r2(A)

- are true, then r8(A, C) is true. In this case, we have r3(Elena, Thomas) and r3(Thomas, Nina) as
- true facts, and we also know that r2(Elena) is true from F5. Therefore, we can apply L5 to get:
- r8(Elena, Nina) is true.
- 1724 Therefore, the answer is True.
- 1725 Statement: r29(Philipp, Elena)
- Answer: We can use logical rule L26: $\forall A, B, C, D : r3(B, A) \land r3(C, B) \land r3(D, C) \land r1(A) \rightarrow r3(B, A) \land r3(B, A) \land r3(B, C) \land r3(B$
- $r^{29}(A, D)$ to infer that if $r^{3}(B, A)$, $r^{3}(C, B)$, $r^{3}(D, C)$, and $r^{1}(A)$ are true, then $r^{29}(A, D)$ is also
- 1728 true
- From the given facts, we know that r^3 (Elena, Thomas) (F46) and r^3 (Thomas, Amelie) (F47).
- Therefore, we can use L26 with A = Philipp, B = Thomas, C = Amelie, and D = Elena to get:
- r3(Thomas, Philipp) (from F62 and r3(B, A))
- r3(Amelie, Thomas) (from F47 and r3(C, B))
- r3(Elena, Amelie) (from F47 and r3(D, C))
- r1(Philipp) (from F14 and r1(A))
- Therefore, we can conclude that r29(Philipp, Elena) is True.
- 1736 Statement: r6(Luisa, Helga)
- Answer: We can use logical rule L3: $\forall A, B : r3(A, B) \land r2(A) \rightarrow r6(A, B)$ to deduce the truth
- value of the statement r6(Luisa, Helga).
- From the facts, we know that r^3 (Luisa, Helga) and r^2 (Luisa) are true. Therefore, we can apply L3 to
- 1740 conclude that r6(Luisa, Helga) is true.
- 1741 Thus, the answer is True.
- 1742 Statement: r14(Amelie, Alina)
- Answer: We can use logical rule L11: $\forall A, B, C, D, E : r3(B, A) \land r3(B, C) \land r3(C, D) \land r3(C, D)$
- $r3(D,E) \wedge r2(A) \rightarrow r14(A,E)$ to infer r14(Amelie, Alina) if we can find facts that satisfy the
- 1745 antecedent of the rule.
- From the facts, we can see that r3 (Amelie, Thomas) and r3 (Thomas, Helga) and r3 (Helga, David) and
- r3(David, Alina) are true. Therefore, we can apply L11 with A =Amelie, B =Thomas, C =Helga,
- 1748 D = David, and E = Alina to get r14(Amelie, Alina) is True.
- Therefore, the answer is True.
- 1750 Statement: r26(Alina, Luisa)
- Answer: We can use logical rule L23: $\forall A, B, C : r3(B, A) \land r3(C, B) \land r2(A) \rightarrow r26(A, C)$ to
- determine the truth value of the statement r26(Alina, Luisa).
- From the facts, we know that r3(Alina, David) and r3(David, Luisa) (F34 and F52). Therefore, we
- can substitute A = Alina, B = David, and C = Luisa into L23 to get:
- 1755 $r3(David, Alina) \wedge r3(Luisa, David) \wedge r2(Alina) \rightarrow r26(Alina, Luisa)$
- r3(Alina, David) is true (F34) and r2(Alina) is also true (F12). Therefore, we can simplify the above
- expression to:
- 1758 $r3(Luisa, David) \rightarrow r26(Alina, Luisa)$
- r3(Luisa, David) is true (F52), so we can conclude that r26(Alina, Luisa) is true.
- 1760 Therefore, the answer is True.

1761 J.2 Few-Shot-CoT examples with Manually-Designed

- 1762 Statement: r14(Amelie, Jonathan)
- Answer: We can use logical rule L11: $\forall A, B, C, D, E : r3(B, A) \land r3(B, C) \land r3(C, D) \land r3(C, D)$
- $r3(D,E) \wedge r2(A) \rightarrow r14(A,E)$. We can substitute A = Amelie and E = Jonathan
- to get: $r3(B,Amelie) \wedge r3(B,C) \wedge r3(C,D) \wedge r3(D,Jonathan) \wedge r2(Amelie) \rightarrow$
- r14(Amelie, Jonathan). Now, we need to find all facts that match the conditions r3(B, Amelie),
- r3(B, C), r3(C, D), r3(D, Jonathan) and r2(Amelie). We have:
- 1768 F41: *r*3(Fabian, Amelie)
- 1769 F47: *r*3(Elena, Amelie)
- 1770 F56: r3(Patrick, Jonathan)
- 1771 F61: *r*3(Emilia, Jonathan)
- 1772 F7: r2(Amelie)
- We can see that F41 and F47 match the condition r3(B, Amelie) with B being Fabian and Elena,
- respectively. F56 and F61 match the condition r3(D, Jonathan) with D being Patrick and Emilia, re-
- spectively. F7 matches the condition r2(Amelie). Therefore, we have two possible ways of assigning

Table 6: The deductive reasoning results of one Symbolic Tree KG with ChatGPT. Results are in %.

	Auto-Randomly	Auto-Filter	Manually-Designed
KG_1	55.6	53.6	40.7

- values to B, C, and D:
- 1777 (1) B = Fabian, D = Patrick.
- 1778 (2) B = Elena, D = Emilia.
- 1779 (3) B = Fabian, D = Emilia.
- (4) B = Elena, D = Patrick.
- Using the values of B and D from the possible assignments above, we can find the value of C that
- satisfies the conditions r3(B, C) and r3(C, D). For assignment (1), we have:
- 1783 F40: *r*3(Fabian, Thomas)
- 1784 F50: *r*3(Thomas, Patrick)
- Since all the conditions are satisfied, we can conclude that the statement r14(Amelie, Jonathan) is
- True if we assign B = Fabian and C = Thomas and D = Patrick. Thus, the output is True.
- 1787 Statement: r31(Philipp, Nina)
- Answer: We can use logical rule L28: $\forall A, B, C, D : r3(B, A) \land r3(C, B) \land r3(C, D) \land r1(A) \rightarrow r3(C, B) \land r3(C, D) \land r3$
- r31(A,D). We can substitute A = Philipp and D = Nina to get: $r3(B,Philipp) \wedge r3(C,B) \wedge r3(D,B)$
- $r3(C,Nina) \wedge r1(Philipp) \rightarrow r31(Philipp,Nina)$. Now, we need to find all facts that match the
- conditions r3(B, Philipp), r3(C, B), r3(C, Nina), and r1(Philipp). We have:
- 1792 F57: *r*3(Patrick, Philipp)
- 1793 F62: r3(Emilia, Philipp)
- 1794 F49: r3(Thomas, Nina)
- 1795 F52: r3(Luisa, Nina)
- We can see that F57 and F62 match the condition r3(B, Philipp) with B being Patrick and Emilia,
- respectively. F49 and F52 match the condition r3(C, Nina) with C being Thomas and Luisa, respec-
- tively. Therefore, there is only one single way of assigning values to B and C.
- 1799 B = Patrick and C = Thomas. We can substitute B and C to get: $r3(Patrick, Philipp) \land$
- 1800 $r3(Thomas, Patrick) \wedge r3(Thomas, Nina) \wedge r1(Philipp) \rightarrow r31(Philipp, Nina)$. Since all
- the conditions are satisfied, we can conclude that the statement r31 (Philipp, Nina) is True if we
- assign B = Patrick and C = Thomas. Thus, the output is True.
- 1803 Statement: r8(Leonie, Nico)
- 1804 Answer: We can use logical rule L5: $\forall A, B, C : r3(A, B) \land r3(B, C) \land r2(A) \rightarrow r8(A, C)$. We
- can substitute A = Leonie and C = Nico to get: $r3(Leonie, B) \land r3(B, Nico) \land r2(Leonie) \rightarrow r3(B, Nico) \land r2(Leonie) \rightarrow r3(B, Nico) \land r3(B, Nic$
- r8(Leonie, Nico). Now, we need to find all facts that match the conditions r3(Leonie, B), r3(B, Conie, B)
- Nico), and r2(Leonie). We have:
- 1808 F43: *r*3(Leonie, Emily)
- 1809 F38: r3(Emily, Nico)
- 1810 F39: r3(Konstantin, Nico)
- 1811 F24: r2(Leonie)
- We can see that F43 matches the condition r3(Leonie, B) with B being Emily. and F38 matches the
- condition r3(B, Nico) with B being Emily. Therefore, we can assign B = Emily. Also, F24 matches
- the condition r2 (Leonie). Therefore, we can conclude that the statement r8 (Leonie, Nico) is True if
- we assign B = Emily. Therefore, the output is True.
- 1816 Statement: r7(Patrick, Alina)
- Answer: We can use logical rule L4: $\forall A, B : r3(A, B) \land r1(A) \rightarrow r7(A, B)$. We can substitute A
- 1818 = Patrick and B = Alina to get: $r3(Patrick, Alina) \wedge r1(Patrick) \rightarrow r7(Patrick, Alina)$. Now,
- we need to find all facts that match the conditions r3(Patrick, Alina) and r1(Patrick). We have:
- 1820 F55: *r*3(Patrick, Alina)
- 1821 F9: *r*1(Patrick)
- We can see that F55 matches the condition r3(Patrick, Alina), and F9 matches the condition
- r_{1823} r_{1} (Patrick). Therefore, we can conclude that the statement r_{1} (Patrick, Alina) is True based on
- 1824 L4. Thus, the output is True.

1825 K Reasoning results of after human selection

In this section, we perform all experiments after human selection, providing the LLMs with only the relevant facts related to the predicted fact.

1828 For example, we have facts:

1843

1863

```
F1: Lea is female.
1829
     F2: Philipp is male.
1830
1831
     F3: Leonie is female.
     F4: Patrick is parent of Alina.
1832
     F5: Patrick is parent of Philipp.
     F6: Nico is parent of David.
1834
     F7: Nico is parent of Lea.
1835
1836
     F8: Laura is parent of Fabian.
1837
     F9: Laura is parent of Felix.
```

If we predict the unknown fact "Philipp is the brother of Alina", we filter out the irrelevant facts and retain only the related facts:

```
1840 F1: Philipp is male.
1841 F2: Patrick is parent of Alina.
1842 F3: Patrick is parent of Philipp.
```

K.1 Comparison of Removing rules and Semantics

To explore the ability of LLMs to utilize internal and external knowledge, we conduct an additional experiment where we provided LLMs with only the relevant facts related to the predicted fact. We compared the performance of *Removing rules* (leveraging internal knowledge) with *Semantics* (providing external logical rules). Surprisingly, from Table 7. we find that *Removing rules* performed better than *Semantics*. This suggests that LLMs possess the necessary internal knowledge to support answering questions and reasoning tasks, and leveraging internal knowledge is more effective for reasoning than relying on external logical rules.

Table 7: Comparison of Removing and providing rules of deductive reasoning. Results are in %.

	Remove rules	Provide rules	
Zero-Shot	80.5	77.1	

1851 LLMs' output examples are provided below:

Semantics: According to the given logical rule, for A to be a granddaughter of D, B must be the parent of A and C must be the parent of B. However, there is no direct connection between Emilia and Lea in the given facts. Although Emilia is the parent of Alina, who is the parent of Lea, Emilia is not the parent of Lea's parent. Therefore, we cannot conclude that Lea is Emilia's granddaughter based on the given facts and logical rule.

Removing rules: Based on the given facts, Emilia is the parent of Alina, who is the parent of Lea.
Therefore, Emilia is the grandmother of Lea, making Lea Emilia's granddaughter.

From the example, we can observe that when relying on external logical rules, LLMs need to strictly adhere to the reasoning process, which can be more challenging for LLMs to predict unknown answers compared to utilizing the commonsense knowledge already contained within LLMs. This suggests that leveraging the internal knowledge of LLMs can be more effective for reasoning tasks.

K.2 Reasoning results after human selection

We conduct deductive and inductive reasoning experiments to examine the performance of LLMs when only provided with the relevant facts related to the predicted fact. The results are presented in Table 8. They demonstrate that after selecting useful information, LLMs perform reasoning tasks more effectively. This finding suggests that LLMs face challenges when processing excessively long in-context information. Selecting relevant facts helps to reduce the memorization load on LLMs and enables them to focus on the most relevant information for reasoning, leading to improved performance.

Table 8: Reasoning results after removing irrelevant information. Results are %.

		Zero-Shot	Zero-Shot-CoT
Deductive	standard	52.6	56.1
	removing irr	55.7	63.0
Inductive	standard	7.14	7.14
	removing irr	67.9	67.9

Table 9: The reasoning results of Symbolic Tree (ChatGPT). Results are in %.

Category	Baseline	deduction	induction	abduction
Logic language	Zero-Shot Zero-Shot-CoT Few-Shot-CoT	52.6 56.1 53.7	7.14 7.14 -	1.95 3.57 13.3
Natural language	Zero-Shot Zero-Shot-CoT Few-Shot-CoT	50.6 50.2 51.9	3.57 7.14	3.90 1.95 8.13

L Reasoning with natural language

1871

1876

1877

1878

1879

1880

1881

1882

1883

1884

1885

1886

In this section, we conducted experiments using the *Symbols* setting with deduction, induction, and abduction on a Symbolic Tree dataset expressed in natural language. The results are presented in Table 9. We observed that, in general, LLMs performed better when using logical language compared to natural language.

M Reasoning results of two representations

For the Symbolic Tree dataset, facts and rules can be represented as logic language and natural language text as the input of LLMs. For example, the fact "motherOf(Alice, Bob)" can be represented as "Alice is Bob's mother"; the fact "r1(Alice, Bob) can be represented as "Alice is r1 of Bob"; the rule " $\forall x, y:$ parentOf(x, y) \rightarrow childOf(y, x)" can be represented as "If x is parent of y, then y is parent of x.". Through numerous trials, we find that for the *Symbols* or *Counter-CS* setting, LLMs tend to perform better when using logic language representations. Conversely, for the *Semantics* setting, LLMs tend to perform better when using natural language text. The results are presented in Table 10. These observations suggest that natural language representations better stimulate the semantic understanding capabilities of LLMs, while logical language representations are more conducive to symbolic reasoning.

Table 10: Deductive reasoning results in different representations. Results are %.

		Zero-Shot	Zero-Shot-CoT
Symbols	logic	52.6	56.1
	natural language	49.0	51.1
Semantics	logic	61.4	61.9
	natural language	69.3	64.3
Counter-CS	logic	52.6	54.4
	natural language	48.7	48.3

1887 N Reasoning with garbled symbols

In this section, we randomly sample 4-8 letters to construct a garbled symbols word as each relation label. However, because LLMs process text by tokens, common sequences of characters found in the text, this random label is often tokenized as a single letter by LLMs. This leads to difficulty for LLMs to comprehend the word-level meaning of the label, harming the performance of reasoning.

Table 11: The deductive reasoning results of one Symbolic Tree KG with ChatGPT. Results are in %.

	Zero-Shot	Zero-Shot-CoT	Few-Shot-CoT	Zero-Plus-Few-Shot-CoT
garbled symbols	52.4	54.4	55.5	52.9
single-token symbols	52.5	53.2	53.9	52.3
ID symbols	54.5	57.0	57.2	52.3
short rules long rules	52.1	51.1	52.9	54.6
	54.5	57.0	57.2	52.3

- We also attempted to use single-token garbled symbols as relation labels, but similar challenges were encountered, with the labels often being interpreted as typos by LLMs, further hindering their
- reasoning abilities. Our experiments and case studies also demonstrate this point.
- As shown in Table 11, using ID symbols (*e.g.*, "r1") as relation labels performs better than using garbled symbols or single-token symbols.
- 1897 Some examples with garbled symbols:
- Statement: ckxovsaq(Alina, Philipp)
- Answer: We can use logical rule L1 to make a deduction about Alina and Philipp. L1 states that
- if there is a person B who is a friend of both A and C, and A is not outgoing, then A trusts C. We
- know that Alina is not outgoing (F12: ntoea(Alina)) and that Philipp is a friend of Emilia (F62:
- yufevh(Emilia, Philipp)). However, we do not have any information about Alina and Emilia's
- friendship. Therefore, we cannot use L1 to make a deduction about ckxovsaq(Alina, Philipp). The
- 1904 answer is False.
- 1905 Statement: *nyheg*(Emily, Nico)
- 1906 Answer: We can use logical rule L3: $\forall A, B : yufevh(A, B) \land ntoea(A) \rightarrow nyheg(A, B)$ to
- determine if the statement is true or false. This rule states that if A is not outgoing and there is a
- 1908 friendship relation between A and B, then A is jealous of B. From the facts, we know that Emily is not
- outgoing (F17: ntoea(Emily)) and she is friends with Nico (F38: yufevh(Emily, Nico)). Therefore,
- we can apply L3 to conclude that nyheq (Emily, Nico) is true. Therefore, the answer is True.
- 1911 Statement: *phnd*(Luisa, Jonathan)
- 1912 Answer: We can use logical rule L5: $\forall A, B, C : yufevh(A, B) \land yufevh(B, C) \land ntoea(A) \rightarrow$
- phnd(A,C) to determine if the statement is true or false. This rule states that if A is not outgoing
- and there is a path from A to C through B, then A is phoned by C. From the given facts, we know that
- Luisa is not outgoing (F8: ntoea(Luisa)). We also know that there is a path from Luisa to Jonathan
- through Patrick (F53: yufevh(Luisa, Patrick) and F56: yufevh(Patrick, Jonathan)). Therefore, we
- can apply L5 with A = Luisa, B = Patrick, and C = Jonathan to get: phnd(Luisa, Jonathan) is true.
- 1918 Therefore, the answer is True.
- 1919 Some examples with single-token symbols:
- 1920 Statement: *iance*(Patrick, Lea)
- Answer: We can use logical rule L4: $\forall A, B : inely(A, B) \land atis(A) \rightarrow lesai(A, B)$, which states
- that if A is an artist and B is in A's inner circle, then B is a muse of A. We can also use fact F53:
- inely(Luisa, Patrick), which means that Luisa is in Patrick's inner circle. Therefore, we can conclude
- that *lesai*(Patrick, Luisa) is true.
- Next, we can use logical rule L6: $\forall A, B, C : inely(A, B) \land inely(B, C) \land atis(A) \rightarrow iance(A, C)$,
- which states that if A is an artist, B is in A's inner circle, and C is in B's inner circle, then A and C
- have an <u>artistic alliance</u>. We can use fact F35: *inely*(Alina, Lea), which means that Alina is in Lea's inner circle. Therefore, we can conclude that *iance*(Alina, Lea) is true.
- Finally, we can use logical rule L1: $\forall A, B, C : inely(B, A) \land inely(B, C) \land uitka(A) \rightarrow$
- reib(A, C), which states that if A is a writer, B is in A's inner circle, and B is also in C's in-
- ner circle, then A and C have a literary rivalry. We can use fact F54: inely(Patrick, Samuel), which
- means that Samuel is in Patrick's inner circle. We can also use fact F35: inely(Alina, Lea), which
- means that Alina is in Lea's inner circle. Therefore, we can conclude that reib(Patrick, Lea) is true.
- Since Patrick has a literary rivalry with Lea, but not an artistic alliance, we can conclude that the
- statement *iance*(Patrick, Lea) is False.

Table 12: The deductive reasoning results of Symbolic Tree datasets. Results are in %.

Category	Model	Baseline	S_1	S_2	S_3	S_4	S_5	S_6	S_7	S_8	S_9	S_{10}	Avg.
	Random	-	52.4	50.8	51.3	50.2	49.3	49.1	48.1	52.3	48.4	49.0	50.1
		Zero-Shot	52.6	50.6	50.5	49.5	55.2	53.1	50.0	53.4	56.6	54.0	52.6
Symbols	ChatGPT	Zero-Shot-CoT	56.1	57.0	55.4	57.0	54.5	56.1	55.5	56.9	50.0	58.0	55.7
	ChatGri	Few-Shot-CoT	53.7	56.9	55.2	54.4	55.1	52.0	54.0	55.8	56.8	54.5	54.8
		Zero-Plus-Few-Shot-CoT	53.7	53.6	55.4	51.4	54.0	50.9	54.0	54.2	58.4	54.5	54.0
		Zero-Shot	70.0	64.8	70.4	65.8	61.4	63.8	65.8	67.4	63.0	68.9	66.1
Semantics	ChatGPT	Zero-Shot-CoT	66.7	64.8	64.6	64.1	64.4	67.2	66.5	66.7	64.6	65.4	65.5
Semantics	ChaiGri	Few-Shot-CoT	71.8	70.4	63.9	69.2	66.7	59.3	68.7	68.3	67.9	64.4	67.1
		Zero-Plus-Few-Shot-CoT	71.3	67.8	66.6	69.5	65.7	60.9	68.4	68.3	66.5	66.8	67.2
	Logic-based	-	100	100	100	100	100	100	100	100	100	100	100

Table 13: The inductive reasoning results of Symbolic Tree datasets. Results are in %.

Category	Model	Baseline	$ S_1 $	S_2	S_3	S_4	S_5	S_6	S ₇	S ₈	S ₉	S ₁₀	Avg.
	ChatGPT	Zero-Shot	7.14	9.09	3.57	7.14	4.54	14.3	4.54	7.14	3.57	0.0	6.10
Symbols	ChaiGFi	Zero-Shot-CoT	7.14	7.14	3.57	14.3	14.3	7.14	3.57	0.0	14.3	7.14	7.86
Symbols	GPT-4	Zero-Shot	14.3	10.7	10.7	10.7	7.14	7.14	10.7	7.14	7.14	7.14	9.28
•	GF 1-4	Zero-Shot-CoT	21.4	7.14	17.9	7.14	3.57	7.14	7.14	7.14	7.14	3.57	8.93
	ChatGPT	Zero-Shot	25.0	32.1	39.3	39.3	42.9	39.3	35.7	32.1	35.7	42.9	36.4
Semantics	ChaiGFI	Zero-Shot-CoT	25.0	28.6	35.7	28.6	35.7	35.7	28.6	35.7	39.3	28.6	32.2
Semantics	GPT-4	Zero-Shot	53.6	53.6	50.0	53.6	50.0	53.6	50.0	57.1	53.6	50.0	52.5
	GF 1-4	Zero-Shot-CoT	53.6	57.1	53.6	53.6	57.1	53.6	50.0	53.6	57.1	50.0	53.9
	Rule-based	-	64.3	60.7	60.7	46.4	67.9	50.0	64.3	57.1	53.6	46.4	57.1

- 1936 Statement: *icers*(Jonathan, David)
- Answer: We can use logical rule L10: $\forall A, B, C, D : inely(B, A) \land inely(B, C) \land inely(C, D) \land inel$
- $atis(A) \rightarrow icers(A, D)$, which states that if A is an artist and B, C, and D are in a line with B in the
- middle, then A influences D.
- 1940 We can use the following facts: *inely*(Jonathan, Alina) (F55)
- 1941 inely(Alina, David) (F34)
- 1942 From these facts, we can conclude that Jonathan, Alina, and David are in a line with Alina in the
- middle. Also, we know that Alina is an artist. Therefore, we can apply L10 to conclude that Alina
- influences David.
- Since Alina influences David, we can use L4: $\forall A, B : inely(A, B) \land atis(A) \rightarrow lesai(A, B)$, which
- states that if A influences B and A is an artist, then B is a musician. Therefore, we can conclude that
- 1947 David is a musician.
- 1948 Therefore, the statement *icers* (Jonathan, David) is False, since David is not an ice skater.
- 1949 Statement: indr(Florian, Nina)
- Answer: We can use logical rule L24: $\forall A, B : inely(B, A) \land atis(A) \rightarrow indr(A, B)$, which states
- that if A is an artist and B is a fan of A, then B is indirectly a fan of all the artists that A is a fan of.
- 1952 From the given facts, we know that atis(Florian) and inely(Emilia, Nina). Therefore, we can
- conclude that Nina is indirectly a fan of Florian, since Emilia is a fan of Florian and Nina is a fan of
- 1954 Emilia.

1956

1959

1955 Thus, the answer is True.

O More reasoning results of Symbolic Tree

We experiment with 10 sampled trees and report the average results in the main body. In this section, we provide the reasoning results of each sampled Symbolic Tree, presented in Table Tables 12 to 14.

P Abductive reasoning on smaller datasets

- 1960 We use smaller Symbolic Tree datasets to conduct the abductive reasoning experiment, which contains
- about 12 entities and 100 facts. The results are provided in Table 15. We compare Symbols and
- 1962 Semantics and find that the Semantics setting still outperforms the Symbols setting. This reinforces
- the hypothesis that preserving semantics enhances the reasoning capabilities of LLMs.

Table 14: The abductive reasoning results of Symbolic Tree KGs. Results are in %.

Category	Model	Baseline	S_1	S_2	S_3	S_4	S_5	S_6	S ₇	S ₈	S_9	S ₁₀	Avg.
		Zero-Shot	1.95	0.31	1.07	1.52	2.36	1.45	1.06	0.75	3.1	1.39	1.50
C	Ch-4CDT	Zero-Shot-CoT	3.57	4.08	5.00	3.03	3.70	3.77	5.28	7.55	7.78	5.21	4.90
Symbols	ChatGPT	Few-Shot-CoT	13.3	7.70	8.39	7.42	10.8	8.55	10.7	14.3	8.95	7.99	9.81
		Zero-Plus-Few-Shot-CoT	22.7	16.7	15.0	11.5	19.9	12.6	12.7	25.3	15.2	16.3	16.8
		Zero-Shot	1.95	3.14	3.57	1.52	2.69	2.32	3.87	3.02	3.89	3.47	2.94
Semantics	ChatGPT	Zero-Shot-CoT	4.22	5.34	4.64	3.63	2.69	2.90	4.23	1.89	3.11	1.39	3.40
		Few-Shot-CoT	9.90	13.2	10.9	7.42	8.59	0.97	11.3	13.0	11.3	11.1	9.77
		Zero-Plus-Few-Shot-CoT	17.5	25.2	22.1	16.7	16.5	18.0	22.2	27.2	22.6	21.5	20.9
	Rule-based	-	100	100	100	100	100	100	100	100	100	100	100

Additionally, abductive reasoning in a shorter context yielded better performance compared to a longer context. This suggests that the length of the context has an impact on reasoning performance. Shorter contexts make selecting relevant and useful information easier while minimizing the influence

Table 15: The abductive reasoning results of a smaller Symbolic Tree. Results are in %.

Category	Baseline	short context	long context
Symbols	ChatGPT: Zero-Shot-CoT	9.78	3.57
	GPT-4: Zero-Shot-CoT	46.7	32.1
Semantics	ChatGPT: Zero-Shot-CoT	5.43	4.22
	GPT-4: Zero-Shot-CoT	59.8	31.8

Replacing entity labels

1964

1965

1966 1967

1968

1974

1975

1980

1981

1982 1983

1984

In this section, we conducted experiments to investigate the effects of replacing entity names (such 1969 as "Alice") with entity IDs (e.g., "e1") in the context of reasoning tasks. The results are provided in 1970 Table 16. Comparing the performance of replacing relation names with replacing both entity and 1971 relation names, we observe that replacing entity names after replacing relation names had little impact 1972 on the overall performance. 1973

Furthermore, we consider the scenario of only replacing entity names. Compared to the case of not replacing any labels, the results indicate that although replacing entity labels retains some level of semantics, it has a detrimental effect on reasoning performance. Additionally, we observed that the negative impact of decoupling the semantics of relations was more significant than that of decoupling the semantics of entities. These findings indicate a substantial portion of the semantic information is concentrated in the relation names. Table 16: Comparison of replacing entity labels in deductive reasoning experiment (ChatGPT).

Results are in %.

	Zero-Shot	Zero-Shot-CoT
replacing none	69.3	66.1
replacing ent	63.6	58.9
replacing rel	54.5	54.5
replacing ent & rel	57.5	55.6

Multi-short rules R

Besides, a single rule can be equivalent to multiple rules. For example, the rule $\forall x, y, z$: $\operatorname{parentOf}(x,y) \wedge \operatorname{parentOf}(y,z) \wedge \operatorname{gender}(x,\operatorname{female}) \rightarrow \operatorname{GrandmotherOf}(x,z)$ can be represented as $\forall x,y,z$: parentOf $(x,y) \land \text{parentOf}(y,z) \rightarrow \text{GrandparentOf}(x,z), \text{GrandparentOf}(x,z) \land$ gender $(x, \text{female}) \rightarrow \text{GrandmotherOf}(x, z)$. We conduct the experiments with both rule representations and find single-longer rules perform better than multiple-short rules. Results are presented in Table 11.