

LLM-BASED TYPED HYPERRESOLUTION FOR COMMONSENSE REASONING WITH KNOWLEDGE BASES

Anonymous authors

Paper under double-blind review

ABSTRACT

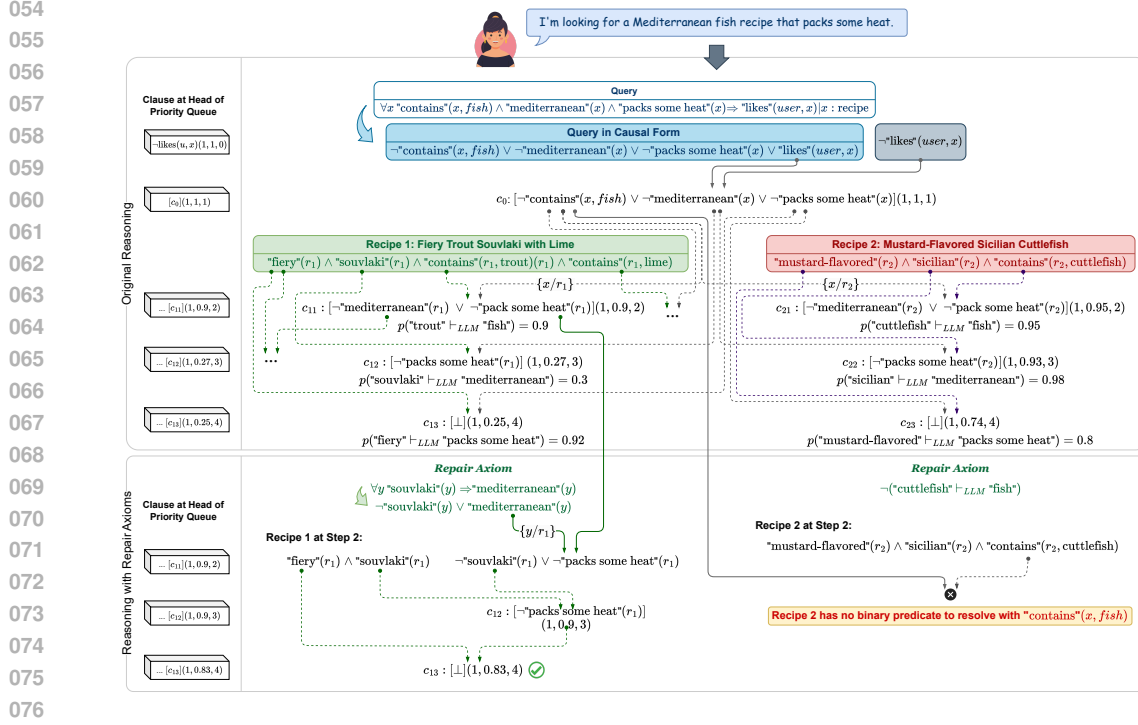
Large language models (LLM) are being increasingly applied to tasks requiring commonsense reasoning. Despite their outstanding potential, the reasoning process of LLMs is prone to errors and hallucinations that hinder their applicability, especially in high-stakes scenarios. Several works have attempted to enhance commonsense reasoning performance of LLMs by (i) using prompting styles that elicit more accurate reasoning, (ii) utilizing the LLM as a semantic parser for a symbolic reasoner, or (iii) enforcing the LLM to simulate a logical inference rule. However, all these solutions have critical limitations: they are unable to leverage the internal commonsense knowledge of the LLM in tandem with an axiomatic knowledge base, they lack a mechanism to reliably repair erroneous inference steps, and their application is restricted to small knowledge bases that fit the context limit of the LLM. In this work, we present LLM-based Typed Hyperresolution (LLM-TH), a logical commonsense reasoning framework that leverages “*theory resolution*”, a concept from classical logical inference which enables integrating LLMs into the “*resolution*” inference rule, thus mitigating reasoning errors and hallucinations and enabling verification of the reasoning procedure. LLM-TH is also equipped with a mechanism for repairing erroneous inference steps supported by theoretical guarantees. Using “*Hyperresolution*” and “*Typed inference*” schemes, we show that LLM-TH can efficiently reason over large knowledge bases consisting of tens of thousands of rules with arbitrary predicate arities. Our experiments on three diverse language-based reasoning tasks—preference reasoning, multi-domain deductive reasoning, and geographical question answering—showcase that LLM-TH, using merely a BART 406M parameter NLI entailment model, significantly reduces reasoning errors compared to baselines using Llama3-70B, Gemini1.5-Flash, GPT-3.5-Turbo, and Mixtral-46.7B.

1 INTRODUCTION

The breakthrough in Large Language Models (LLMs) has significantly impacted AI research, paving the way for deploying AI-powered systems in various tasks and applications. This huge impact is primarily due to the outstanding performance of LLMs in tasks that require substantial reasoning skills (Chang et al., 2024; Plaat et al., 2024). LLMs have also acquired commonsense understanding, a critical element for interacting with the real world (Zhao et al., 2024; Valmeekam et al., 2024). However, reasoning performance of LLMs is not infallible. They commonly show reasoning errors and make hallucinations—generating incorrect outputs that seem valid—which hinders their reliable deployment, particularly in high-risk tasks (Tonmoy et al., 2024; Zhang et al., 2023b).

To overcome these challenges in LLM-based reasoning, several approaches have been proposed in the literature that can be broadly categorized into three groups: (i) Using prompting styles that can elicit more accurate reasoning from the LLM (Wei et al., 2022; Kojima et al., 2022; Zhou et al., 2022) or augmenting the prompt by retrieved information (Lewis et al., 2020b), (ii) using the LLM to translate natural language problem and knowledge bases (KB) for a symbolic logical solver (Olausson et al., 2023; Pan et al., 2023), and (iii) using the LLM to emulate a logical inference rule to solve the reasoning problem (Kazemi et al., 2023; Lee & Hwang, 2024).

These works have notably advanced the logical reasoning performance of LLMs; yet, they are all hindered by a number of important limitations: (a) Their application is limited to small KBs that can



577
578
579
580
581
582
583
584
585

Figure 1: Workflow of LLM-TH shown with a preference reasoning example. Top: Using LLM-based typed hyperresolution to compute proof scores for each recipe option to entail user query. Negated query is the first active clause, and each resolvent is assigned a priority tuple: (type entailment score, predicate entailment score, proof length score) and pushed to the priority queue (only the foremost clause is shown for priority queue of recipe 1 which becomes the next active clause). LLM makes two mistakes: assigning a low score to “Souvlaki” entailing “Mediterranean” and a high score to “cuttlefish” entailing “fish”. Bottom: Both kinds of mistakes can be fixed after the insertion of repair axioms, resulting in the correct recommendation of Recipe 1.

586
587
588
589
590
591
592
593
594

fit into the context limit of the LLM, and are not scalable to perform reasoning on practical KBs such as the widely-used Knowledge Graphs (KG) containing thousands of facts and axioms. (b) They are restricted to perform reasoning on a complete KB containing all rules required to solve the problem. However, assuming access to such KB is typically unrealistic in practical use cases, thus calling for the necessity of a methodology to leverage the internal commonsense knowledge of the LLM in the reasoning process. (c) All steps involved in the reasoning process are not transparent and thus, the correctness of the final answer cannot be determined by inspecting the reasoning process. (d) Upon observation of a reasoning error, they do not provide any reliable framework to fix the error and ensure it will not occur in the future.

595
596

In this work, we aim to address these limitations by making the following contributions:

- 597
598
599
600
601
602
603
604
605
606
607
- We introduce LLM-based Typed Hyperresolution (LLM-TH), a framework for efficient complex logical commonsense reasoning (Singh et al., 2021; Ismayilzade et al., 2023; He et al., 2021), with KBs containing predicates of arbitrary arities, that facilitates the incorporation of the internal commonsense knowledge of LLMs in the reasoning process. LLM-TH is founded on “theory resolution” (Stickel, 1985; Baumgartner, 1992), a concept from classical logical reasoning that allows for the incorporation of specialized theorem provers into the resolution inference rule. (Section 3.1)
 - We equip LLM-TH with a mechanism for incorporating the type information of the variables and constants in the problem domain to prune the proof search space and terminate the exploration of reasoning paths that are unlikely to succeed at very early stages. Also, using hyperresolution, an extension of resolution that enables combining clauses to perform several resolution steps simultaneously, we make LLM-TH an efficient and scalable

108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161

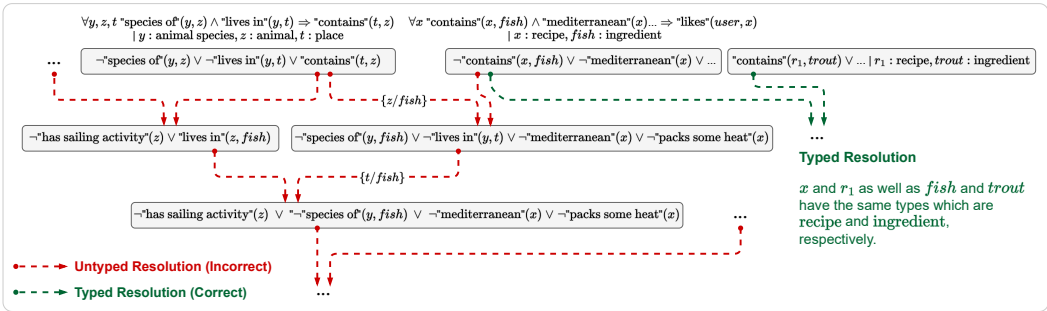


Figure 2: An example of typing mechanism of LLM-TH. The binary predicate “contains” is common among clauses from different domains. Left: Application of resolution rule without considering types leads to resolving a literal about “recipes” with complementary literals from domains such as “animals”. Each incorrect resolution results in a new branch that will be explored but leads to failure, making the process inefficient. Right: in typed resolution, only literals with consistent variable and constant types will be unified, therefore pruning the search space and enhancing efficiency.

reasoning framework for logical commonsense reasoning with LLMs. We show that LLM-TH is easily scalable to KBs consisting of tens of thousands of rules. (Section 5.2.1)

- We show that by providing access to the exact axioms and facts used at every reasoning step, LLM-TH results in a verifiable and faithful reasoning performance. Furthermore, we propose a reliable repair methodology for missed inferences and incorrect reasoning steps due to LLM hallucinations and missed inferences, and provide theoretical proof that it reliably fixes reasoning errors. (Section 4)
- We show that by using the theory hyperresolution framework, LLM-TH is able to leverage the internal commonsense knowledge of the LLM to compensate for KB incompleteness and perform accurate and reliable reasoning. (Section 5.2.2)
- We experiment with LLM-TH on three different tasks involving commonsense reasoning: preference reasoning, multi-domain deductive reasoning, and geographical QA, showcasing the superiority of LLM-TH in terms of answer and reasoning accuracy over Chain of Thought (Kojima et al., 2022; Wei et al., 2022) and retrieval augmented generation (RAG) (Lewis et al., 2020b) baselines using orders of magnitude larger LLMs. (Section 5)

2 RELATED WORKS

Eliciting Stronger Reasoning from LLMs As LLMs scale, they exhibit emergent behaviors such as the capability of solving problems that involve reasoning (Chang et al., 2024; Huang & Chang, 2022). However, their reasoning performance often suffers from errors and incorporating hallucinated facts in their judgments (Tonmoy et al., 2024; Zhang et al., 2023b). Several works have shown that with certain prompting techniques such as Chain of Thought (CoT) (Wei et al., 2022; Kojima et al., 2022), Tree of thought (Yao et al., 2024), Selection-Inference (Creswell et al., 2022), Self-consistency (Wang et al., 2022), Least to most prompting (Zhou et al., 2022), etc., more accurate reasoning can be elicited from LLMs. Retrieval Augmented Generation (RAG) (Lewis et al., 2020b) has also been noted as an effective approach in reducing reasoning errors and hallucinations by including relevant content retrieved from a KB in the prompt to condition LLM’s reasoning on dependable information. These methods have made significant progress in advancing the reasoning behavior of LLMs, but even applying them does not guarantee an accurate reasoning behavior from the LLM. Furthermore, since the LLM is entirely in charge of doing the reasoning, there is no control over the reasoning process and its correctness cannot be verified (Shanahan, 2024; Pan et al., 2023). Moreover, the performance of these methods has been shown to degrade when being applied to long-horizon (Dziri et al., 2024) and out-of-domain (Saparov et al., 2024) reasoning tasks as well as problems involving negation (Anil et al., 2022) and contraposition (Zhang et al., 2024).

Formal Reasoning with LLMs To offer more control over the reasoning process, two groups of work have been proposed for performing formal reasoning with LLMs: (i) *Semantic parsing* meth-

ods remove the LLM from reasoning and only use it to translate the problem into a symbolic format and delegate the reasoning task to a symbolic solver (Pan et al., 2023; Olausson et al., 2023). (ii) Works enforcing the LLM to emulate an established logical inference rule such as backchaining (Kazemi et al., 2023; Lee & Hwang, 2024). While these groups of works make significant progress in mitigating hallucinations, they both rely on the existing rules in the KB and have no particular mechanism to leverage the rich commonsense knowledge of the LLM in their reasoning. Although recent works (Toroghi et al., 2024) have used resolution inference rule for logical LLM-based reasoning, they are restricted to unary predicates and can only do reasoning over small KBs.

3 METHODOLOGY

We first review the (hyper)resolution rule for inference in first-order logic using a function and equality-free first-order logical (FOL) syntax in clausal normal form (Chang & Lee, 2014).¹ We then proceed to introduce the concept of Theory resolution that leverages external sources of knowledge (such as LLMs) to augment the resolution process. Next, we introduce the concept of typed resolution that will be critical for improving the efficiency of inference in our work by prioritizing inference with compatible types and reducing proof length. With all of these ingredients in place, we conclude with our LLM Theory Hyperresolution algorithm (LLM-TH) along with a crucial Lemma that ensures the correctness of *repair axioms* that can be used to override incorrect LLM inferences.

Resolution Rule and Hyperresolution *Resolution* is a sound and complete inference rule which is widely used in logical reasoning. From two premise clauses containing complementary literals, resolution rule derives a *resolvent* clause by canceling (resolving) the complementary literals, e.g.,

$$\frac{A(x) \vee B(x, y) \quad \neg B(w, z) \vee C(z)}{A(x) \vee C(y)}, \quad (1)$$

under the unification $\theta = \{x/w, y/z\}$. Repeated application of the resolution rule will either result in a contradiction, e.g., deriving both $A(x)$ and $\neg A(x)$ indicating an inconsistent clause set, or reaching a point where no further resolutions are possible.

The efficiency of the repeated application of the resolution rule can be substantially improved by *hyperresolution* (Robinson, 1965), an extension of resolution that enables combining multiple resolution steps in one inference step. Concretely, it resolves positive literals with all possible matching negative clauses simultaneously, e.g.,

$$\frac{A(x) \vee B_1(x, y) \dots \vee B_n(x, y) \quad C(z) \vee \neg B_1(w, z) \vee \dots \vee \neg B_n(w, z)}{A(x) \vee C(y)}, \theta = \{x/w, y/z\}. \quad (2)$$

LLM-based Theory Resolution Application of the resolution rule is originally restricted to clauses with complementary literals that share identical predicates. Theory resolution (Stickel, 1985; Baumgartner, 1992) relaxes this condition and broadens the applicability of the resolution rule by integrating special-purpose theories into resolution. Based on theory resolution, given two clauses $c_1 = A(x) \vee B(x, y)$ and $c_2 = \neg C(w, z) \vee D(z)$, if a theorem prover T identifies $B(x, y)$ and $\neg C(w, z)$ under unification $\theta = \{x/w, y/z\}$ to be unsatisfiable (i.e., $\forall x \forall y B(x, y) \wedge \neg C(x, y) \vdash_T \perp$), the clauses can be resolved despite lacking complimentary literals with identical predicates:

$$\frac{A(x) \vee B(x, y) \quad \neg C(w, z) \vee D(z)}{A(x) \vee D(y)}, \theta = \{x/w, y/z\}. \quad (3)$$

In this work, we use an LLM as the theory that identifies the unsatisfiable natural language predicates to perform reasoning via theory resolution. Translating natural language to symbolic form, as semantic parsing methods do, is substantially restricted. For example, they map “*packs some heat*” and being “*spicy*” to completely different symbolic predicates. Therefore, a symbolic reasoner is unable to discern their entailment relationship unless given explicit axioms.

Using LLM-based theory resolution, we can integrate the LLM’s commonsense knowledge into the reasoning process to find entailments between predicates and constants without requiring explicit

¹We refer the reader to Chapters 7-9 of Russell & Norvig(2016) for an in-depth overview of these topics.

axioms. We do this in an extended version of FOL in which predicates, functions, and constants are no longer symbols, but natural language text. In this logical system, which we call *natural language logic*, the unsatisfiability condition in theory resolution reduces to natural language entailment. In other words, if an LLM identifies a natural language predicate B to entail predicate D , i.e., $B(x) \vdash_{LLM} D(x)$, and therefore, $B(x) \wedge \neg D(x) \vdash_{LLM} \perp$, then literals $B(x)$ and $D(x)$ can be resolved. For instance, given clauses $c_1 = \text{“packs some heat”}(x)$ and $c_2 = \neg \text{“spicy”}(x) \vee Q(x)$, in which $Q(x)$ is another literal with a natural language predicate, since the LLM identifies the natural language entailment $\text{“packs some heat”} \vdash_{LLM} \text{“spicy”}$, a theory resolution step can be performed as

$$\frac{\text{“packs some heat”}(x) \quad \neg \text{“spicy”}(x) \vee Q(x)}{Q(x)}. \quad (4)$$

LLM-based Typed Theory Resolution Although resolution is one of the most widely used inference rules in logical reasoning with the key properties of being sound and complete, its application on real-world large-scale knowledge bases is computationally expensive. One of the key challenges is that the space of possible resolutions—the combinations of literals and clauses that can potentially be unified and resolved—can become enormous, and grows exponentially during the resolution process, leading to inefficiencies in finding contradictions or valid derivations. This challenge is often addressed by choosing particular orderings that first explore resolution steps that are more likely to lead to valid proof (Baumgartner, 1992; Sanner & McIlraith, 2006).

In addition to proposing an ordering strategy which is explained in 3.1, we introduce the notion of typing into theory resolution which considerably prunes the space of possible resolutions. Many of the created resolvents during the resolution process, each opening a new search branch, are created by unifying variables and constants with inconsistent ontological types. For instance, a literal $\text{“small”}(x)$ in a clause about *vehicles* can be resolved with a literal $\text{“small”}(y)$ which is about *animals*. However, the search branch created by this resolvent is very unlikely to result in a valid proof as it is created by a semantically invalid unification, and in future steps, the proof will be stopped as no further resolutions could be made. By associating variables and constants with their ontological types and considering type consistency in unification, we can significantly prune the space of allowable resolutions, by preventing the exploration of type inconsistent branches from the beginning.

Variable types can be introduced into an FOL sentence as unary predicates, e.g., $\text{“animal”}(x)$. Considering a universally quantified sentence in FOL as

$$\forall x \forall y H(x) \wedge T(y) \wedge A(x, y) \implies R(x, y), \quad (5)$$

where $H(\cdot)$ and $T(\cdot)$ are unary predicates indicating types of their corresponding variables, we equivalently write the above sentence in the *typed* FOL as

$$\forall x \forall y A(x, y) \implies R(x, y) | x : H, y : T, \quad (6)$$

where, $x : H$ and $y : T$ indicate that x and y are of type H and T respectively. This clause can be equivalently written in the clausal form as

$$\forall x \forall y \neg A(x, y) \vee R(x, y) | x : H, y : T, \quad (7)$$

Next, consider another clause in typed FOL as

$$\forall w \forall z \neg R(w, z) \vee S(w, z) | w : M, z : N. \quad (8)$$

By introducing the unifier $\theta = \{w/x, z/y\}$, we can perform *typed resolution* between the two clauses 7 and 8 as

$$\frac{\neg A(x, y) \vee R(x, y) | x : H, y : T \quad \neg R(w, z) \vee S(w, z) | w : M, z : N}{\neg A(x, y) \vee S(x, y) | x : C, y : D}, \quad (9)$$

where C and D , the types of resolvent variables are $C \equiv H \sqcap M$ and $D \equiv T \sqcap N$ using the following lemma which is proved in Appendix A. Here, \sqcap indicates unary type intersection (conjunction).

Lemma 1. *Resolving two disjunctive clauses c_1 and c_2 that include complimentary literals $l(x_1, \dots, x_n) | x_1 : T_1, \dots, x_n : T_n$ and $\neg l(y_1, \dots, y_n) | y_1 : H_1, \dots, y_n : H_n$ under the unifier $\theta = \{x_1/y_1, \dots, x_2/y_2\}$ creates a resolvent c_{res} with type variables $x_1 : T_1 \sqcap H_1, \dots, x_n : T_n \sqcap H_n$.*

Typed resolution can be directly extended to typed theory resolution in our natural language-based logical system by resolving literals $B(x_1, \dots, x_n) | x_1 : T_1, \dots, x_n : T_n$ and $D(y_1, \dots, y_n) | y_1 : H_1, \dots, H_n$ if the LLM identifies the natural language predicate B to entail D , i.e., $B(x_1, \dots, x_n) \vdash_{LLM} D(y_1, \dots, y_n)$ if their unified variables have consistent types. In the next section, we elaborate on how type consistency is checked in our LLM-TH framework.

Algorithm 1 LLM-TH Algorithm

```

270 1: Input:  $\mathcal{K}, q, \text{max\_proofs}, \text{max\_iters}, \mathcal{I}, \mathcal{F}$ 
271 2:  $\text{proofs} \leftarrow \emptyset$ 
272 3:  $PQ \leftarrow \emptyset$  //  $PQ$  is an initially empty priority queue.
273 4:  $PQ.\text{push}(\neg q, (1, 1, 0))$  // Negation of  $q$  has priority  $(1, 1, 0)$ ,  $PQ$  is ordered by Equation 15
274 5:  $i \leftarrow 0$ 
275 6: while  $(PQ \neq \emptyset) \wedge (|\text{proofs}| \leq \text{max\_proofs}) \wedge (i \leq \text{max\_iters})$  do
276 7:    $c \leftarrow PQ.\text{pop}()$ 
277 8:   if  $c = \perp$  then
278 9:      $\text{proofs} \leftarrow \text{proofs} \cup \{(\mathcal{A}_c, (\rho^t(c), \rho^e(c), \rho^l(c)))\}$  //  $\mathcal{A}_c$  is the set of ancestors of  $c$ , found by
279       backtracking up to  $\neg q$ 
280 10:  else
281 11:     $\beta_c \leftarrow$  candidate clauses in  $\mathcal{K}$  with similar arity and different polarity to  $c$ 
282 12:    for  $c_{\text{target}} \in \beta_c$  do
283 13:      Perform hyperresolution to compute resolvent  $c_{\text{res}}$  of  $c$  and  $c_{\text{target}}$  using Equation 2
284 14:       $PQ.\text{push}(c_{\text{res}}, (\rho^t(c_{\text{res}}), \rho^e(c_{\text{res}}), \rho^l(c_{\text{res}})))$  // cf. Equations 11, 13 and 14
285 15:     $i \leftarrow i + 1$ 
286 16: Output:  $\text{proofs}$ 

```

3.1 LLM-TH ALGORITHM

This section presents LLM-TH, our proposed algorithm for efficient logical commonsense reasoning, which is based on theory resolution (Stickel, 1985; Baumgartner, 1992) extended to LLM theory resolution with predicates of arbitrary arity, hyperresolution, and simultaneous type inference. Its workflow is shown in a worked example in Figures 1 and 2, and formalized in Algorithm 1.

Problem Definition Let \mathcal{Q} denote a set of queries and \mathcal{K} represent a knowledge base (KB) comprising a set of axioms and facts, a set of intended repair axioms \mathcal{I} , and a set of forbidden repair axioms \mathcal{F} , all expressed in natural language logic and clausal form with arbitrary predicate arities. For each query $q \in \mathcal{Q}$, LLM-TH identifies a set of proofs, denoted as proofs . Each proof $f \in \text{proofs}$ consists of a subset of clauses in \mathcal{K} , and is assigned a priority score ρ .

Algorithm To prove \mathcal{K} entails the query q via resolution, we need to show that repeatedly using the resolution rule on $\mathcal{K} \wedge \neg q$ leads to a contradiction, and is thus unsatisfiable. Following the backward chaining paradigm that offers efficiency benefits (Poole & Mackworth, 2010), we pick $\neg q$ as the first active clause to be resolved with a clause from \mathcal{K} . At each step, given an active clause $c = \bigvee_{i=1}^{|c|} l_i$ where each l_i is a literal of arbitrary arity, any clause $c_{\text{target}} \in \mathcal{K}$ as $c_{\text{target}} = \bigvee_{i=1}^{|c_{\text{target}}|} l_{\text{target}_i}$ is considered a *candidate clause* to be theory resolved with c , yielding the resolvent c_{res} , if at least one $(l_i, l_{\text{target}_i})$ pair can be formed where l_i and l_{target_i} have equal arities and different polarities. One can lift this polarity restriction for theory resolutions to allow implicit negation in the language of the predicate (e.g., “not spicy”) in contrast to explicit logical negation (e.g., \neg “spicy”). As \mathcal{K} is often large and expands further with new resolvents being derived as resolution advances, efficiency is a key desideratum, which LLM-TH achieves by prioritizing candidate clauses based on two criteria: (i) type consistency and (ii) predicate entailment between the active and candidate clauses.

Restricting the Space of Resolutions with Typing: The first mechanism used by LLM-TH to improve efficiency is the use of typed theory resolution to restrict the space of allowable resolutions. In typed theory resolution, the types of variables in the resolvent clause are determined by the conjunction of the types of variables in their parent clause, hence LLM-TH prioritizes clauses with variable types that align with the active clause to encourage valid types. For example, if two candidate clauses are considered to be resolved with an active clause of variable type “Animal”, LLM-TH prioritizes a candidate clause with a unifying variable type “Mammal” over one with type “Vehicle”.

LLM-TH leverages the commonsense knowledge of the LLM to obtain the probability of entailment between the variable types. Formally, for the pair of literals with equal arities (l, l_{target}) , denoting the set of argument types of l and l_{target} as $T = \{t_i\}$ and $T' = \{t'_i\}$ respectively, the plausibility of unifying each of their variables can be obtained by calculating the entailment probabilities between

each t_i and t'_i . Since entailment is an asymmetric relation and its direction is unknown, we need to calculate both $t_i \vdash_{LLM} t'_i$ and $t'_i \vdash_{LLM} t_i$ to obtain the type consistency score. The average of type entailment scores for arguments determines $\rho^{type}(c_{res})$, the overall type priority score for c_{res} .

$$\rho^{type}(c_{res}) = \frac{1}{|T_i|} \sum_i (\max\{p(t_i \vdash_{LLM} t'_i), p(t'_i \vdash_{LLM} t_i)\}). \quad (10)$$

Since the main objective is to find the most plausible proofs, i.e., the sequences of most plausible theory resolution steps, we define the first entry of our priority score for each c_{res} as the overall type consistency score of all resolution steps beginning from $\neg q$ that led to its derivation. Let $\mathcal{P}_{c_{res}}$ be the set of parent clauses of c_{res} ; the proof type consistency score of c_{res} is inductively defined as

$$\rho^t(c_{res}) = \left(\prod_{c' \in \mathcal{P}_{c_{res}}} \rho^t(c') \right) \cdot \rho^{type}_{c_{res}}. \quad (11)$$

Resolution Ordering based on Predicate Entailment: Assigning the type consistency scores prunes the resolution search space ensuring that only clauses with compatible argument types will be considered for resolution. To further enhance efficiency, LLM-TH prioritizes the remaining clauses based on their potential for being part of a plausible proof considering their predicate entailment. As explained, in our LLM-based theory resolution framework, LLM entailment is used to identify unsatisfiability of clauses. Therefore, the greater probability the LLM assigns to l_{target} entailing l , the more plausible it becomes to resolve l and l_{target} . Therefore we define the plausibility of a theory resolution step between c and c_{target} by resolving literals l and l_{target} generating c_{res} , denoted by $\rho_{c_{res}}^{entail}$

$$\rho_{c_{res}}^{entail} = \begin{cases} 1 & \text{if } (l_{target} \vdash l) \in \mathcal{I}, \\ 0 & \text{if } \neg(l_{target} \vdash l) \in \mathcal{F}, \\ p(l_{target} \vdash_{LLM} l) & \text{otherwise.} \end{cases} \quad (12)$$

These plausibility scores can help us prioritize the resolvent clauses based on their predicate entailment. For example, in the scenario depicted in Figure 1, resolving “*Sicilian*” with “*Mediterranean*” results in a higher entailment score compared to resolving “*Mustard-flavored*” with “*Mediterranean*”. Hence, it is prudent to prioritize the former resolvent, as it is more likely to contribute to the final proof. As the definition of the overall proof type consistency score, we can compute the overall predicate entailment score inductively to obtain the second entry of our priority score as

$$\rho^e(c_{res}) = \left(\prod_{c' \in \mathcal{P}_{c_{res}}} \rho^e(c') \right) \cdot \rho_{c_{res}}^{entail}. \quad (13)$$

Ultimately, among equally plausible proofs, we prioritize shorter ones that circumvent unnecessary steps. We define the proof length score, a third priority score that accounts for this preference. The proof length score of c_{res} is derived inductively from the maximum length of its parent clauses as

$$\rho^l(c_{res}) = 1 + \max_{c' \in \mathcal{P}_{c_{res}}} \rho^l(c'). \quad (14)$$

Each resolvent c_{res} is assigned the priority tuple $(\rho^t(c_{res}), \rho^e(c_{res}), \rho^l(c_{res}))$ and then pushed to the priority queue PQ , in which the order of clauses is specified as

$$c_1 \preceq c_2 \iff [\rho^t(c_1) > \rho^t(c_2)] \vee [(\rho^t(c_1) = \rho^t(c_2)) \wedge (\rho^e(c_1) > \rho^e(c_2))] \vee [(\rho^t(c_1) = \rho^t(c_2)) \wedge (\rho^e(c_1) = \rho^e(c_2)) \wedge (\rho^l(c_1) < \rho^l(c_2))]. \quad (15)$$

By applying this prioritization scheme, the type consistency priority score first applies a hard filter to avoid exploration of resolvents with invalid types, and the predicate entailment and length priorities together enable an efficient inference via LLM-based theory resolution. These efficiency enhancements enable LLM-TH to be applied to large-scale KBs. Furthermore, it is able to reason over incomplete KBs by benefiting from the commonsense knowledge of the LLM to fill in the missing axioms by identifying entailing predicates in the theory resolution process.

At the start of each iteration of LLM-TH, the clause with the highest rank in PQ becomes the active clause. When a resolution step yields a contradiction, the proof and its respective priority score are added to the set of found proofs by backtracking the ancestor clauses. The algorithm

terminates when either a certain number of proofs are found or the maximum number of iterations is exceeded. Notably, LLM-TH is not limited to proving a single query; instead, it finds a set of proofs and assigns a strength score to each. This feature enables it to evaluate the likelihood of each query being entailed, which is critical for tasks requiring ranking, such as answering multiple-choice questions. Furthermore, LLM-TH can reason on incomplete KBs by using the LLM’s commonsense knowledge to fill in the missing axioms by identifying entailed predicates via theory resolution.

4 FIXING ERRONEOUS RESOLUTIONS IN LLM-TH

LLM-TH enables the verification of the reasoning process by providing access to atomic resolution steps. Therefore, if an incorrect theory resolution step is performed due to an erroneous entailment probability assigned by the LLM, the source of failure can be easily identified. Here, we show that by adding a repair rule, LLM-TH can recover from its mistake and fix the reasoning. In the example provided in Figure 1, the LLM’s mistake in assigning a low entailment score for “*Souvlaki*” to entail “*Mediterranean*” leads to a missed inference. Also, incorrectly considering “*cuttlefish*” to entail being a “*fish*” leads to incorrectly resolving these literals. However, introducing the correct axiom $\forall y \text{“Souvlaki”}(y) \implies \text{“Mediterranean”}(y)$ to \mathcal{I} and introducing $\neg(\text{“cuttlefish”} \vdash_{LLM} \text{“fish”})$ to \mathcal{F} fix these mistakes. We formalize this property here and prove it in Appendix B.

Proposition 1. Consider proof P_c^ϕ using axiom ϕ that derives clause c . For any incorrect LLM reasoning axiom ϕ , a Repair Axiom ϕ' can be inserted such that $P_c^{\phi'}$ will be produced before P_c^ϕ .

5 EXPERIMENTS

We empirically evaluate LLM-TH² on three different tasks representing commonsense reasoning with KBs on different datasets and compare it against variations of four different baselines to compare them from different aspects by answering the following questions:

- **RQ1-Scalability:** How does the reasoning performance of LLM-TH compare to baselines when reasoning with large, but complete knowledge bases?
- **RQ2-Reasoning with incomplete KBs:** How effectively do LLM-TH and the baselines use the LLM’s commonsense knowledge to compensate for the incompleteness of the KB?
- **RQ3-Efficiency:** How is the efficiency of LLM-TH influenced by typed hyperresolution?

5.1 TASKS AND DATASETS DESCRIPTION

- **Preference reasoning:** An exemplar commonsense reasoning task is providing recommendations using natural language statements of user preferences and restrictions. For this task, we use Recipe-MPR (Zhang et al., 2023a), a dataset consisting of 500 queries, e.g., “*I want French food, but I’m on a budget*” and five-way recipe options. This dataset covers a broad range of commonsense reasoning skills such as temporal and analogical reasoning.
- **Multi-domain Deductive reasoning:** Since established datasets for logical commonsense reasoning with LLMs, e.g., ProntoQA (Saparov & He, 2022) and COPA-SSE (Brassard et al., 2022), have small KBs with less than 20 facts and axioms per query, we find them insufficient for evaluating the reasoning capability on large KBs. Thus, following the approach in Saparov & He (2022), we create a deductive reasoning dataset using manually written commonsense axioms and ground facts sampled from Wikidata knowledge graph (Vrandečić & Krötzsch, 2014). This dataset contains more than 32k rules about five different domains: Biological entities, foods, vehicles, drugs and diseases, and sports, and 1000 queries that answering them requires 2 to 7 reasoning steps. We release this dataset to encourage research on LLM-based commonsense reasoning on large-scale KBs.
- **Geographical QA:** Using the same approach for generating the multi-domain deductive reasoning dataset, we create a KB about geographical entities, e.g., cities, deserts, museums, etc. containing more than 12k rules and 500 queries which we also release.

²<https://anonymous.4open.science/r/typed-logic-release-476D/>

Table 1: Reasoning performance of methods across the three datasets on complete KBs. The preference reasoning dataset lacks an explicit rule base, making RAG-based baselines irrelevant.

Method	Preference Reasoning			Deductive Reasoning			Geographical QA		
	Accuracy	RS Macro	RS Micro	Accuracy	RS Macro	RS Micro	Accuracy	RS Macro	RS Micro
GPT-3.5 Turbo									
Zero-Shot CoT	0.86±0.04	0.60	0.80	0.69±0.02	0.45	0.48	0.71±0.03	0.60	0.77
Few-Shot CoT	0.87±0.02	0.65	0.81	0.65±0.03	0.45	0.48	0.82±0.02	0.50	0.53
RAG + Zero-Shot CoT	NA	NA	NA	0.68±0.02	0.85	0.92	0.74±0.01	0.80	0.88
RAG + Few-Shot CoT	NA	NA	NA	0.69±0.02	0.65	0.84	0.83±0.02	0.75	0.84
Gemini-1.5-Flash									
Zero-Shot CoT	0.84±0.04	0.60	0.75	0.60±0.01	0.40	0.81	0.78±0.03	0.20	0.51
Few-Shot CoT	0.86±0.02	0.55	0.79	0.66±0.05	0.20	0.71	0.79±0.02	0.25	0.54
RAG + Zero-Shot CoT	NA	NA	NA	0.78±0.02	0.85	0.93	0.79±0.03	0.50	0.72
RAG + Few-Shot CoT	NA	NA	NA	0.86±0.04	0.45	0.72	0.78±0.03	0.25	0.47
Llama3 70B									
Zero-Shot CoT	0.87±0.01	0.55	0.80	0.80±0.03	0.15	0.77	0.78±0.01	0.25	0.57
Few-Shot CoT	0.91±0.01	0.70	0.84	0.78±0.02	0.55	0.58	0.87±0.02	0.45	0.45
RAG + Zero-Shot CoT	NA	NA	NA	0.78±0.01	0.50	0.81	0.87±0.030	0.40	0.65
RAG + Few-Shot CoT	NA	NA	NA	0.80±0.02	0.75	0.80	0.91±0.02	0.65	0.71
Mixtral 46.7B									
Zero-Shot CoT	0.79±0.03	0.60	0.84	0.59±0.02	0.30	0.66	0.71±0.02	0.50	0.70
Few-Shot CoT	0.74±0.02	0.65	0.83	0.67±0.01	0.45	0.53	0.80±0.01	0.45	0.49
RAG + Zero-Shot CoT	NA	NA	NA	0.65±0.02	0.65	0.81	0.66±0.07	0.25	0.51
RAG + Few-Shot CoT	NA	NA	NA	0.46±0.03	0.30	0.43	0.70±0.06	0.50	0.65
VERA (T5 5B)	0.86	NA	NA	0.76	NA	NA	0.68	NA	NA
LLM-TH (BART 406M)	0.84	0.90	0.94	1.00	1.00	1.00	1.00	1.00	1.00

5.2 BASELINES AND EVALUATION

Existing formal reasoning methods with LLMs, i.e., semantic parsing methods and methods emulating inference rules, suffer from two limitations that make them inapplicable to our studied datasets: (i) they cannot leverage the internal commonsense knowledge of the LLM and only rely on an explicit and complete rule base, which Recipe-MPR lacks. (ii) They are limited to small KBs that fit in the LLM context size, but our studied KBs are much larger. We use established methodologies for eliciting more faithful reasoning from the LLMs as our comparison baselines: (a) zero-shot CoT (Kojima et al., 2022), (b) few-shot CoT (Wei et al., 2022), and (c) RAG (Lewis et al., 2020b) using a dense retriever (Song et al., 2020) to find relevant rules from the KB and prompting the LLM with zero-shot CoT and (d) few-shot CoT. We also compare against (e) VERA (Liu et al., 2023), a method for estimating plausibility of commonsense statements. VERA estimates a score in the range of (0, 1), and we consider a score greater (smaller) than 0.5 as a True (False) prediction. Except VERA which uses T5 (Raffel et al., 2020), we use a number of common LLMs for other baselines: (1) Gemini 1.5-flash, (2) Llama3 (70B), (3) Mixtral (56.7B), and (4) GPT3.5 Turbo, while using BART large (Lewis et al., 2020a)³ (406 M) and RoBERTa large (Liu et al., 2019)⁴ (365 M) tuned on MNLI (Williams et al., 2018) dataset to obtain entailment probabilities for LLM-TH. We also use pyDatalog⁵ for performing hyperresolution for grounding on the KB facts and use Gemini 1.5-flash to convert natural language axioms to clausal natural language logic form.

We evaluate the reasoning performance based on both, (1) the correctness of the final answer, measured by *accuracy*, and (2) the correctness of the reasoning process by evaluating proofs using the *reasoning score* (RS) (Kazemi et al., 2023) metric which is manually calculated for 20 randomly chosen responses in which the final answer was correct. RS is typically assessed as a binary decision based on whether the predicted proof is entirely aligned with the ground truth proof (Kazemi et al., 2023; Lee & Hwang, 2024), which leads to both a single erroneous step and wholly flawed reasoning being penalized equally. To achieve a more granular evaluation of the proofs, we calculate both the conventional *macro RS* and following the idea of Min et al. (2023), we propose a metric which we call *micro RS*. Let P be a provided proof and P^* a valid ground truth proof. Using the indicator function \mathbb{I} , we define the micro RS for each query as $RS_{Micro} = \frac{1}{|P|} \sum_{p \in P} \mathbb{I}(p \in P^*)$.

³<https://huggingface.co/facebook/bart-large-mnli>

⁴<https://huggingface.co/FacebookAI/roberta-large-mnli>

⁵<https://pypi.org/project/pyDatalog/>

Table 2: Reasoning performance on incomplete KBs. Numbers in parenthesis indicate the difference with the method’s performance on a complete KB (Table 1) with ↓ (↑) showing a decrease (increase).

Method	Deductive Reasoning		Geographical QA	
	Zero-Shot CoT	Few-Shot CoT	Zero-Shot CoT	Few-Shot CoT
GPT-3.5-Turbo	0.65±0.01 (0.03↓)	0.54±0.04 (0.15↓)	0.60±0.02 (0.14↓)	0.64±0.03 (0.19↓)
Gemini-1.5-Flash	0.73±0.02 (0.05↓)	0.73±0.01 (0.13↓)	0.68±0.02 (0.11↓)	0.66±0.02 (0.11↓)
Llama3 70B	0.77±0.02 (0.01↓)	0.77±0.01 (0.03↓)	0.66±0.02 (0.21↓)	0.66±0.02 (0.25↓)
Mixtral 46.7B	0.50±0.04 (0.14↓)	0.48±0.02 (0.03↑)	0.53±0.03 (0.12↓)	0.58±0.03 (0.13↓)
LLM-TH (using BART 406M)	0.97	-	0.95	-
LLM-TH (using RoBERTa 356M)	0.96	-	0.92	-

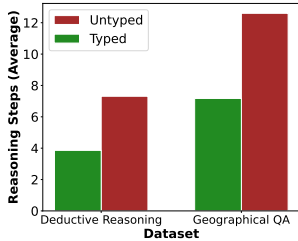


Figure 3: Influence of typing on the efficiency of the inference algorithm.

5.2.1 RQ1: Reasoning Performance on Complete KB Results of the reasoning performance of different methods on complete KBs are provided in Table 1. On the preference reasoning task, which has a small KB, although zero-shot and few-shot CoT with a large LLM such as Llama 3 (70B) yield superior accuracy, LLM-TH outperforms zero-shot and few-shot CoT with Mixtral and is competitive with zero-shot CoT using GPT3.5 and Gemini despite using a much smaller LLM. For the reasoning scores, LLM-TH exhibits a more correct and faithful reasoning process than all other methods. On this task, since the dataset does not contain an explicit KB, RAG-based baselines reduce to zero-shot and few-shot CoT. On larger KBs of deductive reasoning and geographical QA tasks, the limitations of existing LLM-based methods are revealed as none of them compare to the accuracy of LLM-TH. Furthermore, they all obtain imperfect reasoning scores, reflecting their hallucinations and reasoning errors. On these datasets, LLM-TH performs standard resolution which is a sound and complete inference rule, as reflected in the results. While complete KBs are impractical in real-world use cases, results of this experiment highlight that existing baselines, as opposed to LLM-TH, exhibit reasoning failures on large scale KBs *even* when they are complete.

5.2.2 RQ2: Reasoning Performance on Incomplete KBs To enable reasoning over practical KBs, leveraging the commonsense reasoning ability of the LLM to compensate for KB incompleteness is essential. To assess this capability, in our experiments on deductive reasoning and geographical QA datasets that have explicit KBs, we simulate KB incompleteness by omitting one of the rules used in the proof of each query, to test whether the LLM can use its commonsense knowledge to deduce, e.g., “*paying taxes*” implies “*earning revenue*”. Since few-shot and zero-shot CoT methods rely solely on the internal LLM knowledge, they are irrelevant to this RQ that examines the role of KB incompleteness. Hence, we compare LLM-TH against variations of RAG with zero-shot and few-shot CoT prompting. Results in Table 2 show that using the theory hyperresolution framework, LLM-TH is able to achieve significantly higher accuracies compared to the RAG-based baselines that clearly struggle with incompleteness compared to RAG results for complete KBs in Table 1.

5.2.3 RQ3: Influence of Typing on Efficiency To verify the efficiency enhancement offered by introducing type information to the hyperresolution framework of LLM-TH, we perform an ablation experiment on the complete KB, by comparing the average number of reasoning steps that LLM-TH takes to find the answers with its untyped variant that does not consider types in prioritization. In summary, the results of this experiment shown in Figure 3 indicate that typed hyperresolution reduces the proof length by half, which translates to a significantly reduced search space.

6 CONCLUSION

We proposed LLM-TH for logical commonsense reasoning with large and incomplete KBs. Using theory resolution, LLM-TH integrates LLM commonsense knowledge into the resolution inference rule to enable reasoning over incomplete KBs with arbitrary predicates. LLM-TH shows strong performance: it matches or outperforms baselines that use orders of magnitude larger LLMs; its use of an LLM-based typed hyperresolution approach yields high efficiency gains; and its transparency and reparability establish it as a solution for factual and correct reasoning on large-scale KBs. In summary, LLM-TH holds promise to significantly reduce hallucinations in LLM-based reasoning.

Ethics Statement: By introducing LLM-TH, we tried to enhance the transparency and increasing control over the reasoning process of LLM-based logical commonsense reasoning. However, drawing logically valid conclusions does not necessarily mean that all axioms, rules, and the internal commonsense knowledge of the LLM which are leveraged in the reasoning process follow ethical requirements. A responsible and credible usage of LLM-TH, like any other reasoning framework, requires careful considerations and assessments of the knowledge base, the underlying LLM, and the user-defined axioms to ensure desired unbiased and ethical performance.

Reproducibility Statement: We release all our code and data in the supplementary materials, also accessible on the LLM-TH anonymous repository⁶. We also explain the experimental setup and dataset descriptions in Section 5, and include all prompts utilized for the LLM usage in Appendix E, as well as in the supplementary materials.

REFERENCES

- Cem Anil, Yuhuai Wu, Anders Andreassen, Aitor Lewkowycz, Vedant Misra, Vinay Ramasesh, Ambrose Slone, Guy Gur-Ari, Ethan Dyer, and Behnam Neyshabur. Exploring length generalization in large language models. *Advances in Neural Information Processing Systems*, 35:38546–38556, 2022.
- Marco Baroni, Armand Joulin, Allan Jabri, German Kruszewski, Angeliki Lazaridou, Klemen Simoncic, and Tomas Mikolov. Commat: Evaluating the first steps towards a useful general ai. *arXiv preprint arXiv:1701.08954*, 2017.
- Peter Baumgartner. An ordered theory resolution calculus. In *Logic Programming and Automated Reasoning: International Conference LPAR’92 St. Petersburg, Russia, July 15–20, 1992 Proceedings 3*, pp. 119–130. Springer, 1992.
- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 7432–7439, 2020.
- Ana Brassard, Benjamin Heinzerling, Pride Kavumba, and Kentaro Inui. Copa-sse: Semi-structured explanations for commonsense reasoning. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pp. 3994–4000, 2022.
- Chin-Liang Chang and Richard Char-Tung Lee. *Symbolic logic and mechanical theorem proving*. Academic press, 2014.
- Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. A survey on evaluation of large language models. *ACM Transactions on Intelligent Systems and Technology*, 15(3):1–45, 2024.
- Antonia Creswell, Murray Shanahan, and Irina Higgins. Selection-inference: Exploiting large language models for interpretable logical reasoning. *arXiv preprint arXiv:2205.09712*, 2022.
- Ernest Davis and Gary Marcus. Commonsense reasoning and commonsense knowledge in artificial intelligence. *Communications of the ACM*, 58(9):92–103, 2015.
- Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine Li, Liwei Jiang, Bill Yuchen Lin, Sean Welleck, Peter West, Chandra Bhagavatula, Ronan Le Bras, et al. Faith and fate: Limits of transformers on compositionality. *Advances in Neural Information Processing Systems*, 36, 2024.
- Tianqing Fang, Zeming Chen, Yangqiu Song, and Antoine Bosselut. Complex reasoning over logical queries on commonsense knowledge graphs. *arXiv preprint arXiv:2403.07398*, 2024.
- Weinan He, Canming Huang, Yongmei Liu, and Xiaodan Zhu. Winologic: A zero-shot logic-based diagnostic dataset for winograd schema challenge. In *Proceedings of the 2021 conference on empirical methods in natural language processing*, pp. 3779–3789, 2021.

⁶<https://anonymous.4open.science/r/typed-logic-release-476D/>

- 594 Jie Huang and Kevin Chen-Chuan Chang. Towards reasoning in large language models: A survey.
595 *arXiv preprint arXiv:2212.10403*, 2022.
596
- 597 Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong
598 Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hallucination in large language
599 models: Principles, taxonomy, challenges, and open questions. *arXiv preprint arXiv:2311.05232*,
600 2023.
- 601 Mete Ismayilzada, Debjit Paul, Syrielle Montariol, Mor Geva, and Antoine Bosselut. Crow: Bench-
602 marking commonsense reasoning in real-world tasks. *arXiv preprint arXiv:2310.15239*, 2023.
603
- 604 Mehran Kazemi, Najoung Kim, Deepti Bhatia, Xin Xu, and Deepak Ramachandran. Lambada:
605 Backward chaining for automated reasoning in natural language. In *Proceedings of the 61st*
606 *Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp.
607 6547–6568, 2023.
- 608 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large
609 language models are zero-shot reasoners. *Advances in neural information processing systems*,
610 35:22199–22213, 2022.
- 611 Stefanie Krause and Frieder Stolzenburg. Commonsense reasoning and explainable artificial in-
612 telligence using large language models. In *European Conference on Artificial Intelligence*, pp.
613 302–319. Springer, 2023.
- 614 Jinu Lee and Wonseok Hwang. Symba: Symbolic backward chaining for multi-step natural language
615 reasoning. *arXiv preprint arXiv:2402.12806*, 2024.
616
- 617 Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer
618 Levy, Veselin Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-
619 training for natural language generation, translation, and comprehension. In *Proceedings of the*
620 *58th Annual Meeting of the Association for Computational Linguistics*, pp. 7871–7880, 2020a.
- 621 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal,
622 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented gener-
623 ation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:
624 9459–9474, 2020b.
- 625 Hugo Liu and Push Singh. Conceptnet—a practical commonsense reasoning tool-kit. *BT technology*
626 *journal*, 22(4):211–226, 2004.
627
- 628 Jiacheng Liu, Wenya Wang, Dianzhuo Wang, Noah A Smith, Yejin Choi, and Hannaneh Ha-
629 jishirzi. Vera: A general-purpose plausibility estimation model for commonsense statements.
630 *arXiv preprint arXiv:2305.03695*, 2023.
- 631 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike
632 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining
633 approach. *arXiv preprint arXiv:1907.11692*, 2019.
634
- 635 Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke
636 Zettlemoyer, and Hannaneh Hajishirzi. Factscore: Fine-grained atomic evaluation of factual pre-
637 cision in long form text generation. In *Proceedings of the 2023 Conference on Empirical Methods*
638 *in Natural Language Processing*, pp. 12076–12100, 2023.
- 639 Robert C Moore. *The role of logic in knowledge representation and commonsense reasoning*. SRI
640 International. Artificial Intelligence Center, 1982.
- 641 Theo Olausson, Alex Gu, Ben Lipkin, Cedegao Zhang, Armando Solar-Lezama, Joshua Tenenbaum,
642 and Roger Levy. Linc: A neurosymbolic approach for logical reasoning by combining language
643 models with first-order logic provers. In *Proceedings of the 2023 Conference on Empirical Meth-*
644 *ods in Natural Language Processing*, pp. 5153–5176, 2023.
645
- 646 Liangming Pan, Alon Albalak, Xinyi Wang, and William Wang. Logic-lm: Empowering large lan-
647 guage models with symbolic solvers for faithful logical reasoning. In *Findings of the Association*
for Computational Linguistics: EMNLP 2023, pp. 3806–3824, 2023.

- 648 Mete Ismayilzada Debjit Paul, Syrielle Montariol, Mor GevaO, and Antoine Bosselut. Crow: Bench-
649 marking commonsense reasoning in real-world tasks.
650
- 651 Aske Plaat, Annie Wong, Suzan Verberne, Joost Broekens, Niki van Stein, and Thomas Back. Rea-
652 soning with large language models, a survey. *arXiv preprint arXiv:2407.11511*, 2024.
- 653 David L Poole and Alan K Mackworth. *Artificial Intelligence: foundations of computational agents*.
654 Cambridge University Press, 2010.
655
- 656 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi
657 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text
658 transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- 659 John Alan Robinson. A machine-oriented logic based on the resolution principle. *Journal of the*
660 *ACM (JACM)*, 12(1):23–41, 1965.
661
- 662 Stuart J Russell and Peter Norvig. *Artificial intelligence: a modern approach*. Pearson, 2016.
- 663 Scott Sanner and Sheila A McIlraith. An ordered theory resolution calculus for hybrid reasoning in
664 first-order extensions of description logic. In *KR*, pp. 100–111, 2006.
665
- 666 Maarten Sap, Hannah Rashkin, Derek Chen, Ronan LeBras, and Yejin Choi. Socialliqa: Common-
667 sense reasoning about social interactions. *arXiv preprint arXiv:1904.09728*, 2019.
- 668 Abulhair Saparov and He He. Language models are greedy reasoners: A systematic formal analysis
669 of chain-of-thought. In *The Eleventh International Conference on Learning Representations*,
670 2022.
671
- 672 Abulhair Saparov, Richard Yuanzhe Pang, Vishakh Padmakumar, Nitish Joshi, Mehran Kazemi,
673 Najoung Kim, and He He. Testing the general deductive reasoning capacity of large language
674 models using ood examples. *Advances in Neural Information Processing Systems*, 36, 2024.
- 675 Murray Shanahan. Talking about large language models. *Communications of the ACM*, 67(2):68–79,
676 2024.
677
- 678 Ke Shen and Mayank Kejriwal. An experimental study measuring the generalization of fine-tuned
679 language representation models across commonsense reasoning benchmarks. *Expert Systems*, 40
680 (5):e13243, 2023.
- 681 Shikhar Singh, Nuan Wen, Yu Hou, Pegah Alipoormolabashi, Te-Lin Wu, Xuezhe Ma, and Nanyun
682 Peng. Com2sense: A commonsense reasoning benchmark with complementary sentences. *arXiv*
683 *preprint arXiv:2106.00969*, 2021.
684
- 685 Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. MpNet: Masked and permuted pre-
686 training for language understanding. *Advances in neural information processing systems*, 33:
687 16857–16867, 2020.
- 688 Mark E Stickel. Automated deduction by theory resolution. *Journal of Automated Reasoning*, 1(4):
689 333–355, 1985.
690
- 691 SM Tonmoy, SM Zaman, Vinija Jain, Anku Rani, Vipula Rawte, Aman Chadha, and Amitava Das.
692 A comprehensive survey of hallucination mitigation techniques in large language models. *arXiv*
693 *preprint arXiv:2401.01313*, 2024.
- 694 Armin Toroghi, Willis Guo, Ali Pesaranhader, and Scott Sanner. Verifiable, debuggable, and re-
695 pairable commonsense logical reasoning via LLM-based theory resolution. In *The 2024 Confer-*
696 *ence on Empirical Methods in Natural Language Processing*, 2024.
- 697 Karthik Valmeekam, Matthew Marquez, Alberto Olmo, Sarath Sreedharan, and Subbarao Kamb-
698 hampati. Planbench: An extensible benchmark for evaluating large language models on planning
699 and reasoning about change. *Advances in Neural Information Processing Systems*, 36, 2024.
700
- 701 Denny Vrandečić and Markus Krötzsch. Wikidata: a free collaborative knowledgebase. *Communi-*
cations of the ACM, 57(10):78–85, 2014.

- 702 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha
703 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language
704 models. In *The Eleventh International Conference on Learning Representations*, 2022.
- 705 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
706 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*
707 *neural information processing systems*, 35:24824–24837, 2022.
- 708 Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus for sen-
709 tence understanding through inference. In *Proceedings of the 2018 Conference of the North Amer-
710 ican Chapter of the Association for Computational Linguistics: Human Language Technologies,
711 Volume 1 (Long Papers)*, pp. 1112–1122, 2018.
- 712 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik
713 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Ad-
714 vances in Neural Information Processing Systems*, 36, 2024.
- 715 Haochen Zhang, Anton Korikov, Parsa Farinneya, Mohammad Mahdi Abdollah Pour, Manasa
716 Bharadwaj, Ali Pesaraghader, Xi Yu Huang, Yi Xin Lok, Zhaoqi Wang, Nathan Jones, et al.
717 Recipe-mpr: A test collection for evaluating multi-aspect preference-based natural language re-
718 trieval. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Devel-
719 opment in Information Retrieval*, pp. 2744–2753, 2023a.
- 720 Yanfang Zhang, Yiliu Sun, Yibing Zhan, Dapeng Tao, Dacheng Tao, and Chen Gong. Large lan-
721 guage models as an indirect reasoner: Contrapositive and contradiction for automated reasoning.
722 *arXiv preprint arXiv:2402.03667*, 2024.
- 723 Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao,
724 Yu Zhang, Yulong Chen, et al. Siren’s song in the ai ocean: a survey on hallucination in large
725 language models. *arXiv preprint arXiv:2309.01219*, 2023b.
- 726 Zirui Zhao, Wee Sun Lee, and David Hsu. Large language models as commonsense knowledge for
727 large-scale task planning. *Advances in Neural Information Processing Systems*, 36, 2024.
- 728 Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuur-
729 mans, Claire Cui, Olivier Bousquet, Quoc Le, et al. Least-to-most prompting enables complex
730 reasoning in large language models. *arXiv preprint arXiv:2205.10625*, 2022.

734 A PROOF OF LEMMA 1

735 **Lemma 1.** *Resolving two disjunctive clauses c_1 and c_2 that include complimentary literals*
736 *$l(x_1, \dots, x_n)|x_1 : T_1, \dots, x_n : T_n$ and $\neg l(y_1, \dots, y_n)|y_1 : H_1, \dots, y_n : H_n$ under the unifier*
737 *$\theta = \{x_1/y_1, \dots, x_2/y_2\}$ creates a resolvent c_{res} with type variables $x_1 : T_1 \sqcap H_1, \dots, x_n : T_n \sqcap H_n$.*

738 *Proof.* Assume c_1 to be $A(x_1, \dots, x_n) \vee l(x_1, \dots, x_n)|x_1 : T_1, \dots, x_n : T_n$ and c_2 to be $\neg l(y_1, \dots, y_n) \vee$
739 $B(y_1, \dots, y_n)|y_1 : H_1, \dots, y_n : H_n$. Following 5 and 6, we can rewrite these clauses in implication
740 form FOL by adding type predicates as

$$741 \forall x_1, \dots, \forall x_n \bigwedge_{i=1}^n T_i(x_i) \wedge \neg A(x_1, \dots, x_n) \implies l(x_1, \dots, x_n), \quad (16)$$

$$742 \forall y_1, \dots, \forall y_n \bigwedge_{i=1}^n H_i(y_i) \wedge l(y_1, \dots, y_n) \implies B(y_1, \dots, y_n), \quad (17)$$

743 Converting these clauses to the disjunctive form yields

$$744 \forall x_1, \dots, \forall x_n \bigvee_{i=1}^n \neg T_i(x_i) \vee A(x_1, \dots, x_n) \vee l(x_1, \dots, x_n), \quad (18)$$

$$745 \forall y_1, \dots, \forall y_n \bigvee_{i=1}^n \neg H_i(y_i) \vee \neg l(y_1, \dots, y_n) \vee B(y_1, \dots, y_n). \quad (19)$$

Now, we can resolve $l(x_1, \dots, x_n)$ with $\neg l(y_1, \dots, y_n)$ under the unifier $\theta = \{x_1/y_1, \dots, x_n/y_n\}$ as

$$\frac{\bigvee_{i=1}^n \neg T_i(x_i) \vee A(x_1, \dots, x_n) \vee l(x_1, \dots, x_n) \quad \bigvee_{i=1}^n \neg H_i(x_i) \vee \neg l(y_1, \dots, y_n) \vee B(y_1, \dots, y_n)}{\bigvee_{i=1}^n \neg T_i(x_i) \vee \bigvee_{i=1}^n \neg H_i(x_i) \vee A(x_1, \dots, x_n) \vee B(x_1, \dots, x_n)}, \quad (20)$$

which can be rewritten as

$$\forall x_1, \dots, \forall x_n \bigwedge_{i=1}^n T_i(x_i) \wedge H_i(x_i) \implies A(x_1, \dots, x_n) \vee B(x_1, \dots, x_n). \quad (21)$$

Therefore, the unary type predicates for each x_i becomes the conjunction of the types from their parent clauses, which in our typed FOL notation, can be equivalently written as

$$\forall x_1, \dots, \forall x_n A(x_1, \dots, x_n) \vee B(x_1, \dots, x_n) | x_1 : T_1 \sqcap H_1, \dots, x_n : T_n \sqcap H_n. \quad (22)$$

□

B PROOF OF REPAIRABILITY OF LLM-TH

Proposition 1. Consider proof P_c^ϕ using axiom ϕ that derives clause c . For any incorrect LLM reasoning axiom ϕ , a Repair Axiom ϕ' can be inserted such that $P_c^{\phi'}$ will be produced before P_c^ϕ .

Proof. A proof $P_c^\phi = P_c \cup \{\phi\}$ can be viewed as the combined set of clauses P_c and ϕ that derive clause c . An incorrect reasoning reasoning step ϕ can be either a missed inference due to the LLM mistakenly assigning a low priority to a resolution, or an incorrect resolution assigned a high priority due to an LLM misbelief.

We begin the proof for missed inference case, for which the repair axiom $\phi' \in \mathcal{I}$ is available. We can obtain the proof score $\rho^e(P_c^\phi)$ of clause c by inductively unrolling Equation 13 for $\rho^e(c)$ over all ancestor clauses P_c^ϕ that derive it. This yields a simple product form: $\rho^e(P_c^\phi) = \rho_\phi^{entail}$.

$\prod_{c' \in P_c} \rho_{c'}^{entail}$. Now, comparing two different derivations P_c^ϕ and $P_c^{\phi'}$ of c , we can easily show that $\rho^e(P_c^{\phi'}) > \rho^e(P_c^\phi)$ since $\frac{\rho^e(P_c^{\phi'})}{\rho^e(P_c^\phi)} = \frac{\rho_{\phi'}^{entail} \cdot \prod_{c' \in P_c} \rho_{c'}^{entail}}{\rho_\phi^{entail} \cdot \prod_{c' \in P_c} \rho_{c'}^{entail}} = \frac{\rho_{\phi'}^{entail}}{\rho_\phi^{entail}} > 1$ given that the explicit Repair Axiom has $\rho_{\phi'}^{entail} = 1$ (following Equation 13) while the LLM entailment score $\rho_\phi^{entail} < 1$. Hence, the proof $P_c^{\phi'}$ containing the Repair Axiom ϕ' will always be given precedence over P_c^ϕ according to the total ordering of Equation 15 used to prioritize proofs in the LLM-TH Algorithm 1.

Correctness of the proposition for case ϕ is an incorrect resolution that $\phi' \in \mathcal{F}$ forbids it is obvious as in presence of ϕ' , the priority score assigned to the resolvent clause c , i.e., ρ_ϕ^{entail} will be set to 0 (following Equation 13). As a result, the resolvent clause is excluded from any proof by being assigned the lowest priority. □

C EXPLANATION OF THE LLM-TH ALGORITHM

In this section, we provide a detailed explanation of the LLM-TH algorithm, presented in Algorithm 1.

The algorithm takes as inputs the knowledge base (\mathcal{K}), the query (q), a computational budget defined by the maximum number of iterations (max_iters), the maximum number of found proofs shown by max_proofs , and two optional sets of repair axioms: (i) the set of intended resolutions denoted by \mathcal{I} , and (ii) the set of forbidden resolutions shown by \mathcal{F} .

As explained in Section 3.1, LLM-TH is not limited to finding a single proof to the query, but rather it can find a number of proofs with different plausibilities, a capability that is necessary for answering multiple-choice questions. For example, in the preference reasoning example provided in Figure 1, it finds a proof for each of the candidate recipes with a plausibility score assigned to each proof, which is used to rank the recipe options.

The algorithm initializes an empty set, $proofs$, to store discovered proofs, an iteration counter i , and an initially empty priority queue PQ . PQ stores resolvent clauses and orders them according to their priority scores, using Equation 15. In the backward chaining approach used in LLM-TH, the first clause stored in PQ is $\neg q$ with the priority of $(1, 1, 0)$. The algorithm’s main loop executes as long as the priority queue (PQ) is not empty, the maximum proof limit is not reached, and the iteration cap is not exceeded. During each iteration, the clause c with the highest priority in PQ is dequeued and designated as the active clause. In the initial iteration, since PQ contains only $\neg q$, this clause becomes the first active clause.

If the active clause is \perp , the previous resolution step has resulted in a contradiction. The proof can then be reconstructed by tracing the ancestors of c , starting from its parent set \mathcal{P}_c and continuing up to the initial clause $\neg q$. Once a proof is found, it is added to the set $proofs$. However, if $c \neq \perp$, the algorithm proceeds by selecting clauses to resolve with c . It first identifies β_c , the set of candidate clauses in \mathcal{K} that are resolvable with c . A clause must meet two conditions to be resolvable with c in a theory resolution step following Equation 3: (i) it must share the same arity as c , and (ii) it must have a different polarity from c . After forming β_c , LLM-TH iterates through all clauses $c_{target} \in \beta_c$ and resolves each with c via hyperresolution, as described in Equation 2, to obtain c_{res} . It then forms the priority tuple using Equations 11, 13, and 14, and adds c_{res} to PQ based on this priority.

The algorithm continues until either the maximum number of proofs is found, PQ becomes empty, or the maximum allowed number of iterations is reached.

D COMMONSENSE REASONING

The general conception that humans have about the world and how it works, often referred to as commonsense knowledge, and the ability to reason about it, known as commonsense reasoning, are essential capabilities that average human beings possess. In order to be able to interact with the real world and humans, AI agents are also required to obtain this knowledge (Baroni et al., 2017; Shen & Kejriwal, 2023). Therefore, commonsense reasoning has been a central problem in AI over the past decades (Moore, 1982; Liu & Singh, 2004; Davis & Marcus, 2015). With the advent of language models, and particularly, transformer-based architectures, models were developed that exhibited a decent performance on some commonsense reasoning tasks such as physical and social reasoning (Bisk et al., 2020; Sap et al., 2019). Further improvements were obtained by scaled language models, i.e., LLMs, as they exhibited an outstanding capability of commonsense understanding and reasoning across various evaluation benchmarks (Zhao et al., 2024; Krause & Stolzenburg, 2023).

These improvements paved the way for inspecting new aspects of commonsense reasoning, such as reasoning on problems involving logical inference, a vein of work referred to as *complex commonsense reasoning* (Fang et al., 2024; He et al., 2021). The introduction of this line of work shed light on important limitations of LLMs in challenging commonsense reasoning tasks. Although LLMs have shown outstanding reasoning performance on simple commonsense reasoning benchmarks, they struggle with more challenging tasks, such as commonsense reasoning in real-world use cases (Paul et al.).

864 These remaining challenges and the risk of reasoning errors restrict the deployability of LLMs in
 865 real-world applications, especially in high-stakes scenarios. Despite the extensive amount of exist-
 866 ing work on commonsense reasoning with LLMs, to the best of our knowledge, no work has targeted
 867 the problem of combining commonsense reasoning with reasoning over factual KBs. This problem
 868 is particularly important since LLMs show factual hallucinations when they lack domain-specific
 869 knowledge(Huang et al., 2023; Tonmoy et al., 2024). In this work, we aim to take complex log-
 870 ical commonsense reasoning task a step further, by proposing a methodology that leverages both
 871 commonsense reasoning knowledge of the LLM, as well as the factual knowledge of the KB.

873 E PROMPTS USED FOR LANGUAGE MODELS

874 We provide the prompts that we used for the LLMs in the experiments of this paper. They are also
 875 included in our repository along with the implementation code and data.

876 Prompt for Preference Reasoning Task (Zero-shot and Few-shot CoT)

877 Task: You will be given a query that asks for a recipe and five
 878 options that you have to choose from. Think step by step. First
 879 state your reason for your choice and then say: "Therefore, the
 880 selected recipe is <recipe id>.
 881 Query: {{QUERY}}
 882 [Examples if few-shot:]

883 Prompt for Multi-domain Deductive Reasoning Task (Zero-shot and Few-shot CoT)

884 Task: You will be given a query about some knowledge graph
 885 entities in the form of a first order logic predicate that is
 886 either True or False (for example, "CanHoldIn(Apple, Basket)"
 887 which means one can hold an apple in a basket). Your task is to
 888 identify whether the answer to this query is "True" or "False"
 889 and also provide a proof of the answer. First, state your proof
 890 mentioning the rules you used and then say: "Therefore, the
 891 answer is True" or "Therefore, the answer is False". Think step
 892 by step.
 893 Query: {{QUERY}}
 894 [Examples if few-shot:]

895 Prompt for Multi-domain Deductive Reasoning (RAG with Zero-shot and Few-shot CoT)

896 Task: You will be given a query about some knowledge graph
 897 entities in the form of a first order logic predicate that is
 898 either True or False (for example, "CanHoldIn(Apple, Basket)"
 899 which means one can hold an apple in a basket) and a Knowledge
 900 Base containing a set of rules that will help you identify the
 901 answer. Your task is to identify whether the answer to the query
 902 is "True" or "False" and also provide a proof of the query using
 903 the knowledge base. First state your proof mentioning the rules
 904 you used and then say: "Therefore, the answer is True" or
 905 "Therefore, the answer is False". Think step by step.
 906 Query: {{QUERY}}
 907 KB: {{KB}}
 908 [Examples if few-shot:]

909 Prompt for Geographical QA Task (Zero-shot and Few-shot CoT)

910 Task: You will be given a query about geographical entities in
 911 the form of a first order logic predicate that is either True or
 912 False. Your task is to identify whether the answer to the query
 913 is "True" or "False" and also provide a proof of the query. First
 914 state your proof mentioning the rules you used and then say:
 915 "Therefore, the answer is True" or "Therefore, the answer is
 916 False". Think step by step.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

```
Query: {{QUERY}}  
[Examples if few-shot:]
```

Prompt for Geographical QA Task (RAG with Zero-shot and Few-shot CoT)

```
Task: You will be given a query about geographical entities in  
the form of a first order logic predicate that is either True or  
False, and a Knowledge Base containing a set of rules that will  
help you identify the answer. Your task is to identify whether  
the answer to the query is "True" or "False" and also provide a  
proof of the query using the knowledge base. First state your  
proof mentioning the rules you used and then say: "Therefore, the  
answer is True" or "Therefore, the answer is False". Think step  
by step.  
Query: {{QUERY}}  
KB: {{KB}}  
[Examples if few-shot:]
```

F ANECDOTAL EXAMPLES

To offer more insight into the responses and proofs provided by LLM-TH and the comparison baselines, this section presents anecdotal examples illustrating each model’s performance on the evaluated tasks. Since Llama3-70B shows the best overall performance for our comparison baselines in the experimental results provided in Section 5, we show the exemplar outputs from this LLM. Specifically, we showcase the outputs from the following models:

- LLM-TRes
- Zero-shot Chain of Thought prompting
- Few-shot Chain of Thought prompting
- RAG + Zero-shot Chain of Thought prompting
- RAG + Few-shot Chain of Thought prompting

We apply these models to the Deductive reasoning and Geographical QA tasks studied in RQ1 and RQ2 and provide anecdotal examples to elucidate their capabilities and limitations. In each anecdotal example, we present the query, the set of relevant facts and axioms from the KB, and the ground truth proof, followed by the response that each model provides. We also indicate correct proof steps with green color and highlight incorrect ones in red. It is noteworthy that since RQ2 requires access to the KB and studies the model’s ability to compensate for KB incompleteness, the only relevant baselines are RAG-based models.

Examples of the experiments on the complete KB are provided in section F.1, and examples of performance on the incomplete KB are detailed in section F.2.

F.1 EXPERIMENTS ON THE COMPLETE KB

F.1.1 Deductive Reasoning

Query: Does Sulfadiazine disinfect Rhodospirillaceae?

Relevant Rules from the KB:

Facts:

- (1) “Subclass of”(“Rhodospirillaceae” “Rhodospirillaceae bacteria”),
- (2) “Gram Negative”(“Rhodospirillaceae bacteria”)
- (3) “Treats Gram Negative”(“Sulfadiazine”)

Axioms:

- (4)

$$\forall A, B, C, \text{“Subclass of”}(A, B) \wedge \text{“Gram Negative”}(B) \wedge \text{“Treats Gram Negative”}(C) \\ \implies \text{“Unresistant To”}(A, C) | A: \text{“Bacteria”}, B: \text{“Taxon”}, C: \text{“Antibiotic”}$$
- (5)

$$\forall A, B, \text{“Unresistant To”}(A, B) \implies \text{“Can be Used to Treat”}(A, B) | A: \text{“Bacteria”}, \\ B: \text{“Antibiotic”}$$
- (6)

$$\forall A, B, \text{“Can be Used to Treat”}(A, B) \implies \text{“Can be Used to Kill”}(A, B) | A: \text{“Bacteria”}, \\ B: \text{“Antibiotic”}$$

(7) $\forall A, B, \text{“Can be Used to Kill”}(A, B) \implies \text{“Disinfects”}(A, B) \mid A: \text{“Bacteria”}, B: \text{“Antibiotic”}$

Ground Truth Proof:

1- Applying (1), (2), (3) to (4) with the following substitution:

$\theta = \{A/\text{“Rhodospirillaceae”}, B/\text{“Rhodospirillaceae bacteria”}, C/\text{“Sulfadiazine”}\}$ yields:

“Unresistant To”(“Rhodospirillaceae”, “Sulfadiazine”).

2- Applying “Unresistant To”(“Rhodospirillaceae”, “Sulfadiazine”) to (5)

with $\theta = \{A/\text{“Rhodospirillaceae”}, B/\text{“Sulfadiazine”}\}$ yields:

“Can be Used to Treat”(“Rhodospirillaceae”, “Sulfadiazine”)

3-Applying “Can be Used to Treat”(“Rhodospirillaceae”, “Sulfadiazine”) to (6)

with $\theta = \{A/\text{“Rhodospirillaceae”}, B/\text{“Sulfadiazine”}\}$ yields:

“Can be Used to Kill”(“Rhodospirillaceae”, “Sulfadiazine”)

4- Applying “Can be Used to Kill”(“Rhodospirillaceae”, “Sulfadiazine”) to (7)

with $\theta = \{A/\text{“Rhodospirillaceae”}, B/\text{“Sulfadiazine”}\}$ yields:

5- “Disinfects”(“Rhodospirillaceae”, “Sulfadiazine”).

Therefore, the answer to the query is True.

LLM-TH:

query: “Disinfects”(“Rhodospirillaceae”, “Sulfadiazine”):

Negated query: \neg “Disinfects”(“Rhodospirillaceae”, “Sulfadiazine”)

Note: In the following sequent calculus formulae, we use “R” to refer to “Rhodospirillaceae”, “S” to represent “Sulfadiazine”, “GN” to show “Gram Negative”, “TGN”, to denote “Treats Gram Negative”, and “SO” to show “Subclass Of” due to space limitations.

Step 1-

$$\frac{\neg \text{“Disinfects”}(“R”, “S”) \quad \neg \text{“Can be Used to Kill”}(A, B) \vee \text{“Disinfects”}(A, B)}{\neg \text{“Can be Used to Kill”}(“R”, “S”)}$$

Under unification $\theta = \{A/\text{“Rhodospirillaceae”}, B/\text{“Sulfadiazine”}\}$.

Step 2-

$$\frac{\neg \text{“Can be Used to Kill”}(“R”, “S”) \quad \text{“Can be Used to Kill”}(A, B) \vee \neg \text{“Can be Used to Treat”}(A, B)}{\neg \text{“Can be Used to Treat”}(“R”, “S”)}$$

Under unification $\theta = \{A/\text{“Rhodospirillaceae”}, B/\text{“Sulfadiazine”}\}$.

Step 3-

$$\frac{\neg \text{“Can be Used to Treat”}(“R”, “S”) \quad \text{“Can be Used to Treat”}(A, B) \vee \neg \text{“Unresistant To”}(A, B)}{\neg \text{“Unresistant To”}(“R”, “S”)}$$

Step 4-

$$\frac{\neg \text{“Unresistant To”}(“R”, “S”) \quad \text{“Unresistant To”}(A, B) \vee \neg \text{“SO”}(A, B) \vee \neg \text{“GN”}(B) \vee \neg \text{“TGN”}(C)}{\neg \text{“SO”}(“R”, B) \vee \neg \text{“GN”}(B) \vee \neg \text{“TGN”}(“S”)}$$

Under unification $\theta = \{A/\text{“Rhodospirillaceae”}, C/\text{“Sulfadiazine”}\}$.

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

Step 5-

$$\frac{\neg "SO"("R", B) \vee \neg "GN"(B) \vee \neg "TGN"("S") \quad "SO"("R", "R \textit{ bacteria}") \wedge "GN"("R \textit{ bacteria"}) \wedge "TGN"("S")}{\perp}$$

Proof score: $\rho_{\text{query}} = (1, 1, 5)$

* * *

Zero-shot CoT:

To prove the answer to the query, I will use the following rules from the knowledge base:

1. Rhodospirillaceae is a type of Bacteria. (IsA(Rhodospirillaceae, Bacteria))
 2. Sulfadiazine is a type of Antibiotic. (IsA(sulfadiazine, Antibiotic))
 3. Antibiotics disinfect Bacteria. (Disinfects(Antibiotic, Bacteria))
- Using modus ponens on rules 1, 2, and 3, I can infer that: Disinfects(sulfadiazine, Rhodospirillaceae)
Therefore, the answer is **True**.

Few-shot CoT:

1. Sulfadiazine is an antibacterial drug.
2. Rhodospirillaceae is a family of bacteria.
3. Sulfadiazine is not typically used as a disinfectant, but rather as an antibiotic to treat bacterial infections.

Therefore, the answer is **False**.

RAG + Zero-shot CoT:

1. SubClassOf(Rhodospirillaceae, Rhodospirillaceae bacteria) (KB)
2. GramNegative(Rhodospirillaceae bacteria) (KB)
3. $\forall A:\textit{bacteria}, B:\textit{Taxon}, C:\textit{antibiotic}, \textit{SubClassOf}(A, B) \ \&\& \ \textit{GramNegative}(B) \ \&\& \ \textit{TreatsGramNegative}(C) \implies \textit{UnResistantTo}(A, C)$ (KB)
4. UnResistantTo(Rhodospirillaceae, sulfadiazine) (inferred from 1, 2, and 3, with A = Rhodospirillaceae, B = Rhodospirillaceae bacteria, and C = sulfadiazine)
5. $\forall A:\textit{bacteria} \ B:\textit{antibiotic}, \textit{UnResistantTo}(A, B) \implies \textit{CanBeUsedToTreat}(A, B)$ (KB)
6. CanBeUsedToTreat(Rhodospirillaceae, sulfadiazine) (inferred from 4 and 5)
7. $\forall A:\textit{bacteria} \ B:\textit{antibiotic}, \textit{CanBeUsedToTreat}(A, B) \implies \textit{CanBeUsedToKill}(A, B)$ (KB)
8. CanBeUsedToKill(Rhodospirillaceae, sulfadiazine) (inferred from 6 and 7)
9. $\forall A:\textit{bacteria} \ B:\textit{antibiotic}, \textit{CanBeUsedToKill}(A, B) \implies \textit{Disinfects}(A, B)$ (KB)
10. Disinfects(Rhodospirillaceae, sulfadiazine) (inferred from 8 and 9)

Therefore, the answer is **True**.

RAG + Few-shot CoT:

1. Rhodospirillaceae is a gram negative bacteria,
2. sulfadiazine treats gram negative bacteria,
3. and sulfadiazine can be used to treat gram negative bacteria, therefore sulfadiazine can be used to treat Rhodospirillaceae.

Therefore, the answer is **True**.

Discussion:

1134 These anecdotal examples demonstrate that on a complete knowledge base (KB), LLM-TH can
 1135 successfully prove the query by resolving complementary clauses, without requiring the LLM to
 1136 generate any additional axioms.
 1137 An analysis of the proofs generated by the baseline methods reveals limitations in LLM reasoning,
 1138 particularly in zero-shot and few-shot settings. For instance, the zero-shot CoT baseline asserts that
 1139 “Antibiotics disinfect bacteria”; however, this statement is overly broad, as not all antibiotics can
 1140 disinfect all types of bacteria. While the final answer in this case is correct, there is a flaw in one of
 1141 the reasoning steps. Also, the few-shot CoT baseline produces an incorrect final answer due to an
 1142 error in reasoning. In contrast, the RAG-based baselines, utilizing both prompting styles, correctly
 1143 prove the query by leveraging retrieved facts from the KB.

1144 F.1.2 Geographical QA

1145
 1146 **Query:** Can I drive between Biyabanak and Ryam Frolovskiy?

1147
 1148 **Relevant Rules from the KB:**

1149 Facts:

- 1150
 1151 (1) “*Located in*”(“*Biyabanak*”, “*Iran*”),
 1152 (2) “*Located in*”(“*Ryam Frolovskiy*”, “*Russia*”),
 1153 (3) “*Share Land Borders*”(“*Iran*”, “*Azerbaijan*”)
 1154 (4) “*Share Land Borders*”(“*Azerbaijan*”, “*Russia*”)
 1155

1156 Axioms:

- 1157 (5)
 1158 $\forall A, B, \text{“Share Land Borders”}(A, B) \implies \text{“Land Connected”}(A, B) | A : \text{country}, B : \text{country}$
 1159
 1160 (6)
 1161 $\forall A, B, C, \text{“Land Connected”}(A, B) \wedge \text{“Land Connected”}(B, C)$
 1162 $\implies \text{“Land Connected”}(A, C) | A : \text{location}, B : \text{location}, C : \text{location}$
 1163
 1164 (7)
 1165 $\forall A, B, C, D, \text{“Located in”}(A, C) \wedge \text{“Located in”}(B, D) \wedge \text{“Land Connected”}(C, D)$
 1166 $\implies \text{“Can Drive Between”}(A, B) | A : \text{location}, B : \text{location}, C : \text{location}, D : \text{location}$
 1167
 1168

1169 **Ground Truth Proof:**

1170 1- Applying (3) to (5) with the following substitution:

1171 $\theta = \{A/\text{“Iran”}, B/\text{“Azerbaijan”}\}$ yields:

1172 (8) “*Land Connected*”(“*Iran*”, “*Azerbaijan*”).

1173 2- Applying (3) to (5) with the following substitution:

1174 $\theta = \{A/\text{“Azerbaijan”}, B/\text{“Russia”}\}$ yields:

1175 (9) “*Land Connected*”(“*Azerbaijan*”, “*Russia*”).

1176 3- Applying (8), (9) to (6) with the following substitution:

1177 $\theta = \{A/\text{“Iran”}, B/\text{“Azerbaijan”}, C/\text{“Russia”}\}$ yields:

1178 (10) “*Land Connected*”(“*Iran*”, “*Russia*”).

4- Applying (1), (2), (10) to (7) with the following substitution:

$\theta = \{A/“Biyabanak”, B/“Ryam Frolovskiy”, C/“Iran”, D/“Russia”\}$ yields:

(11) “Can Drive Between”(“Biyabanak”, “Ryam Frolovskiy”).

Therefore, the answer to the query is True.

LLM-TH:

Proof for query “Can Drive Between”(“Biyabanak”, “Ryam Frolovskiy”)

Negated query: \neg “Can Drive Between”(“Biyabanak”, “Ryam Frolovskiy”)

Note: In the following sequent calculus formulae, we use “Bi” to refer to “Biyabanak”, “RF” to represent “Ryam Frolovskiy”, “Az” to refer to ‘Azerbaijan’, “CDB” to show “Can Drive Between”, “LI” to show “Located In”, “LC” to show “Land Connected”, and “SLB” to show “Share Land Borders” due to space limitations.

Step 1-

$$\frac{\neg“CDB”(“Bi”, “RF”) \quad “CDB”(A,B) \vee \neg“LI”(A,C) \vee \neg“LI”(B,D) \vee \neg“LC”(C,D)}{\neg“LI”(“Bi”,C) \vee \neg“LI”(“RF”,D) \vee \neg“LC”(C,D)}$$

Under unification $\theta = \{A/“Biyabanak”, C/“Ryam Frolovskiy”\}$.

Step 2-

$$\frac{\neg“LI”(“Bi”,C) \vee \neg“LI”(“RF”,D) \vee \neg“LC”(C,D) \quad “LI”(“Bi”, “Iran”) \wedge “LI”(“RF”, “Russia”)}{\neg“LC”(“Iran”, “Russia”)}$$

Under unification $\theta = \{C/“Iran”, D/“Russia”\}$.

Step 3-

$$\frac{\neg“LC”(“Iran”, “Russia”) \quad “LC”(A,C) \vee \neg“LC”(A,B) \vee \neg“LC”(B,C)}{\neg“LC”(“Iran”,B) \vee \neg“LC”(B, “Russia”)}$$

Under unification $\theta = \{A/“Iran”, C/“Russia”\}$.

Step 4-

$$\frac{\neg“LC”(“Iran”,B) \quad “LC”(“Iran”,B) \vee \neg“SLB”(“Iran”,B)}{\neg“SLB”(“Iran”,B)}$$

Step 5-

$$\frac{\neg“LC”(B, “Russia”) \quad “LC”(B, “Russia”) \vee \neg“SLB”(B, “Russia”)}{\neg“SLB”(B, “Russia”)}$$

Step 6-

$$\frac{\neg“SLB”(“Iran”, B) \vee \neg“SLB”(B, “Russia”) \quad “SLB”(“Iran”, “Az”) \wedge “SLB”(“Az”, “Russia”)}{\perp}$$

Under unification $\theta = \{B/“Az”\}$.

Proof score: $\rho_{\text{query}} = (1, 1, 6)$

Zero-shot CoT:

To answer this query, I will use the following rules from the knowledge base:

1. CanDriveBetween(x, y) \implies Connected(x, y)

2. Connected(Biyabanak, Ryam Frolovskiy) \implies False (this is a fact in the knowledge base)

Using rule 1, I can infer that `CanDriveBetween(Biyabanak, Ryam Frolovskiy)` implies `Connected(Biyabanak, Ryam Frolovskiy)`.
 However, according to rule 2, `Connected(Biyabanak, Ryam Frolovskiy)` is `False`.
 Therefore, the answer is `False`.

Few-shot CoT:

1. `Biyabanak` is a desert in Iran,
 2. and `Ryam Frolovskiy` is an island in the Caspian Sea, which is not connected to Iran by land.
 Therefore, the answer is `False`.

RAG + Zero-shot CoT:

1. `LocatedIn(Biyabanak, Iran)` (Given in KB)
 2. `PartOf(Iran, Iran)` (Rule: $\forall A:\text{place or location } B:\text{place or location, } \text{LocatedIn}(A, B) \implies \text{PartOf}(A, B)$)
 3. `ShareLandBorders(Iran, Azerbaijan)` (Given in KB)
 4. `LandConnected(Iran, Azerbaijan)` (Rule: $\forall A:\text{country } B:\text{country, } \text{ShareLandBorders}(A, B) \implies \text{LandConnected}(A, B)$)
 5. `LocatedIn(Ryam Frolovskiy, Azerbaijan)` (Assuming this is true, since it's not given in KB)
 6. `PartOf(Ryam Frolovskiy, Azerbaijan)` (Rule: $\forall A:\text{place or location } B:\text{place or location, } \text{LocatedIn}(A, B) \implies \text{PartOf}(A, B)$)
 7. `CanDriveBetween(Biyabanak, Ryam Frolovskiy)` (Rule: $\forall A, B, C, D, \text{LocatedIn}(A, C) \ \&\& \ \text{LocatedIn}(B, D) \ \&\& \ \text{LandConnected}(C, D) \implies \text{CanDriveBetween}(A, B) | A:\text{location, } B:\text{location, } C:\text{location, } D:\text{location}$)
 Therefore, the answer is `True`.

RAG + Few-shot CoT:

1. Since `Biyabanak` is a desert in Iran,
 2. and `Ryam Frolovskiy` is a place in Russia,
 3. and Iran and Russia are land connected (although not explicitly stated in the provided KB, but it can be inferred from the general knowledge that Iran and Russia share a border), then it is possible to drive between `Biyabanak` and `Ryam Frolovskiy`.
 Using the rule: $\forall A, B, C, D, \text{LocatedIn}(A, ?C) \ \&\& \ \text{LocatedIn}(B, ?D) \ \&\& \ \text{LandConnected}(?C, ?D) \implies \text{CanDriveBetween}(A, B)$ With $A = \text{Biyabanak}, B = \text{Ryam Frolovskiy}, C = \text{Iran}, D = \text{Russia}$
 Therefore, the answer is `True`.

Discussion:

Similar to the deductive reasoning experiment, LLM-TH produces the correct proof on a complete knowledge base (KB) without requiring the LLM to surface commonsense axioms from it. This example also highlights the limitations of the LLM-based baselines. For instance, the zero-shot CoT baseline incorrectly claims that *Biyabanak* and *Ryam Frolovskiy* are not connected, leading to an incorrect final answer. Similarly, the few-shot CoT baseline provides an inaccurate statement about *Ryam Frolovskiy*, describing it as an island in the *Caspian sea* when it is, in fact, a forest in *Russia*. The RAG-based baselines further reveal challenges in the retrieval stage. For example, the relevant fact about the location of *Ryam Frolovskiy* is not retrieved, prompting the LLM to make assumptions in the zero-shot setting. Although the final answer in this case is correct, the reasoning steps contain errors and deviate from the correct proof. In the few-shot prompting setting, while the final answer is also correct, it is based on an inaccurate statement (in fact, Iran and Russia do not share a land border).

1296 F.2 EXPERIMENTS ON THE INCOMPLETE KB

1297

1298 **F.2.1 Deductive Reasoning**

1299

1300 **Query:** Does Sulfadiazine resolve an infection caused by Rhodospirillaceae?

1301

1302 **Relevant Rules from the KB:**

1303 Facts:

1304

1305 (1) “Subclass of”(“Rhodospirillaceae” “Rhodospirillaceae bacteria”),

1306 (2) “Gram Negative”(“Rhodospirillaceae bacteria”)

1307

1308 (3) “Treats Gram Negative”(“Sulfadiazine”)

1309 Axioms:

1310

1311 (4)

1312 $\forall A, B, C, \text{“Subclass of”}(A, B) \wedge \text{“Gram Negative”}(B) \wedge \text{“Treats Gram Negative”}(C)$
1313 $\implies \text{“Unresistant To”}(A, C) | A: \text{“Bacteria”}, B: \text{“Taxon”}, C: \text{“Antibiotic”}$

1314 (5)

1315 $\forall A, B, \text{“Unresistant To”}(A, B) \implies \text{“Can be Used to Treat”}(A, B) | A: \text{“Bacteria”},$
1316 “B: Antibiotic”

1317

1318

1319

1320 (6)

1321 $\forall A, B, \text{“Can be Used to Treat”}(A, B) \implies \text{“Can be Used to Kill”}(A, B) | A: \text{“Bacteria”},$
1322 “B: Antibiotic”

1323

1324

1325

1326 (7) $\forall A, B, \text{“Can be Used to Kill”}(A, B) \implies \text{“Disinfects”}(A, B) | A: \text{“Bacteria”}, \text{“B: Antibiotic”}$

1327

1328 **Ground Truth Proof:**

1329

1330 1- Applying (1), (2), (3) to (4) with the following substitution:

1331 $\theta = \{A/\text{“Rhodospirillaceae”}, B/\text{“Rhodospirillaceae bacteria”}, C/\text{“Sulfadiazine”}\}$ yields:

1332

1333 “Unresistant To”(“Rhodospirillaceae”, “Sulfadiazine”).

1334 2- Applying “Unresistant To”(“Rhodospirillaceae”, “Sulfadiazine”) to (5)

1335

1336 with $\theta = \{A/\text{“Rhodospirillaceae”}, B/\text{“Sulfadiazine”}\}$ yields:

1337 “Can be Used to Treat”(“Rhodospirillaceae”, “Sulfadiazine”)

1338

1339 3-Applying “Can be Used to Treat”(“Rhodospirillaceae”, “Sulfadiazine”) to (6)

1340 with $\theta = \{A/\text{“Rhodospirillaceae”}, B/\text{“Sulfadiazine”}\}$ yields:

1341

1342 “Can be Used to Kill”(“Rhodospirillaceae”, “Sulfadiazine”)

1343 4- Applying “Can be Used to Kill”(“Rhodospirillaceae”, “Sulfadiazine”) to (7)

1344

1345 with $\theta = \{A/\text{“Rhodospirillaceae”}, B/\text{“Sulfadiazine”}\}$ yields:

1346 5- “Disinfects”(“Rhodospirillaceae”, “Sulfadiazine”).

1347

1348

1349

1350 6- Using the commonsense axiom (absent from the KB) that:

$$1351 \quad \forall A, B, \text{“Disinfects”}(A, B)$$

$$1352 \quad \implies \text{“Resolves Infection Caused By”}(A, B) | A: \text{“Bacteria”}, B: \text{“Antibiotic”}$$

1353 Therefore, the answer to the query is True.

1356 **LLM-TH:**

1357 Proof for query “Resolves Infection Caused By”(“Rhodospirillaceae”, “Sulfadiazine”):

1358 Negated query: \neg “Resolves Infection Caused By”(“Rhodospirillaceae”, “Sulfadiazine”)

1359 Note: In the following sequent calculus formulae, we use “R” to refer to “Rhodospirillaceae”,
1360 “S” to represent “Sulfadiazine”, “GN” to show “Gram Negative”, “TGN”, to denote
1361 “Treats Gram Negative”, “SO” to show “Subclass Of” due to space limitations.

1362 **Step 1-**

$$1363 \quad \frac{1364 \quad \neg \text{“Resolves Infection Caused By”}(“R”, “S”) \quad \text{“Disinfects}(A, B)” \vee \neg \text{“Can be Used to Kill}(A, B)”}{1365 \quad \neg \text{“Can be Used to Kill”}(“R”, “S”)}$$

1366 Under unification $\theta = \{A/\text{“Rhodospirillaceae”}, B/\text{“Sulfadiazine”}\}$.

1367 In this step, by theory resolving \neg “Resolves Infection Caused By”(“R”, “S”) with
1368 “Disinfects”(A, B), the LLM successfully surfaced the commonsense axiom:

$$1369 \quad \forall A, B, \text{“Disinfects”}(A, B)$$

$$1370 \quad \implies \text{“Resolves Infection Caused By”}(A, B) | A: \text{“Bacteria”}, B: \text{“Antibiotic”}.$$

1371 **Step 2-**

$$1372 \quad \frac{1373 \quad \neg \text{“Can be Used to Kill”}(“R”, “S”) \quad \text{“Can be Used to Kill”}(A, B) \vee \neg \text{“Can be Used to Treat”}(A, B)}{1374 \quad \neg \text{“Can be Used to Treat”}(“R”, “S”)}$$

1375 Under unification $\theta = \{A/\text{“Rhodospirillaceae”}, B/\text{“Sulfadiazine”}\}$.

1376 **Step 3-**

$$1377 \quad \frac{1378 \quad \neg \text{“Can be Used to Treat”}(“R”, “S”) \quad \text{“Can be Used to Treat”}(A, B) \vee \neg \text{“Unresistant To”}(A, B)}{1379 \quad \neg \text{“Unresistant To”}(“R”, “S”)}$$

1380 **Step 4-**

$$1381 \quad \frac{1382 \quad \neg \text{“Unresistant To”}(“R”, “S”) \quad \text{“Unresistant To”}(A, B) \vee \neg \text{“SO”}(A, B) \vee \neg \text{“GN”}(B) \vee \neg \text{“TGN”}(C)}{1383 \quad \neg \text{“SO”}(“R”, B) \vee \neg \text{“GN”}(B) \vee \neg \text{“TGN”}(“S”)}$$

1384 Under unification $\theta = \{A/\text{“Rhodospirillaceae”}, C/\text{“Sulfadiazine”}\}$.

1385 **Step 5-**

$$1386 \quad \frac{1387 \quad \neg \text{“SO”}(“R”, B) \vee \neg \text{“GN”}(B) \vee \neg \text{“TGN”}(“S”) \quad \text{“SO”}(“R”, \text{“R bacteria”}) \wedge \text{“GN”}(“R bacteria”) \wedge \text{“TGN”}(“S”)}{1388 \quad \perp}$$

1389 Therefore, the query is proven to be True with the Proof score: $\rho_{\text{query}} = (1, 0.981, 5)$

1390 * * *

1400 **RAG + Zero-shot CoT:**

1401 1. SubClassOf(Rhodospirillaceae, Rhodospirillaceae bacteria) (from KB)

- 1404 2. $\text{GramNegative}(\text{Rhodospirillaceae bacteria})$ (from KB)
- 1405
- 1406 3. $\forall A, B, C, \text{SubClassOf}(A, B) \wedge \text{GramNegative}(B) \wedge \text{TreatsGramNegative}(C) \implies$
- 1407 $\text{UnResistantTo}(A, C)$ (from KB)
- 1408
- 1409 4. $\text{SubClassOf}(\text{Rhodospirillaceae}, \text{Rhodospirillaceae bacteria}) \implies$
- 1410 $\text{PartOf}(\text{Rhodospirillaceae}, \text{Rhodospirillaceae bacteria})$ (from KB, using the rule
- 1411 $\forall A, B, \text{SubClassOf}(A, B) \implies \text{PartOf}(A, B)|A: \text{animal}, B: \text{Taxon}$)
- 1412
- 1413 5. $\text{PartOf}(\text{Rhodospirillaceae}, \text{Rhodospirillaceae bacteria}) \wedge \text{GramNegative}(\text{Rhodospirillaceae bacteria}) \wedge$
- 1414 $\text{TreatsGramNegative}(\text{sulfadiazine}) \implies \text{UnResistantTo}(\text{Rhodospirillaceae}, \text{sulfadiazine})$ (from 2,
- 1415 4, and the assumption that sulfadiazine treats Gram-Negative bacteria)
- 1416
- 1417 6. $\forall A, B, \text{UnResistantTo}(A, B) \implies \text{CanBeUsedToTreat}(A, B)|A: \text{bacteria}, B: \text{antibiotic}$ (from KB)
- 1418
- 1419 7. $\text{UnResistantTo}(\text{Rhodospirillaceae}, \text{sulfadiazine}) \implies$
- 1420 $\text{CanBeUsedToTreat}(\text{Rhodospirillaceae}, \text{sulfadiazine})$ (from 5 and 6)
- 1421
- 1422 8. $\forall A, B, \text{CanBeUsedToTreat}(A, B) \implies \text{CanBeUsedToKill}(A, B)|A: \text{bacteria}, B: \text{antibiotic}$ (from
- 1423 KB)
- 1424
- 1425 9. $\text{CanBeUsedToTreat}(\text{Rhodospirillaceae}, \text{sulfadiazine}) \implies$
- 1426 $\text{CanBeUsedToKill}(\text{Rhodospirillaceae}, \text{sulfadiazine})$ (from 7 and 8)
- 1427
- 1428 10. $\forall A, B, \text{CanBeUsedToKill}(A, B) \implies \text{Disinfects}(A, B)|A: \text{bacteria}, B: \text{antibiotic}$ (from KB)
- 1429
- 1430 11. $\text{CanBeUsedToKill}(\text{Rhodospirillaceae}, \text{sulfadiazine}) \implies$
- 1431 $\text{Disinfects}(\text{Rhodospirillaceae}, \text{sulfadiazine})$ (from 9 and 10) Therefore, the answer is True.

RAG + Few-shot CoT:

- 1428
- 1429 1- Rhodospirillaceae is a gram-negative bacteria,
- 1430
- 1431 2- and sulfadiazine treats gram-negative bacteria. Therefore, the answer is **True**.
-

Discussion:

1433

1434 This example demonstrates how LLM-TH can perform commonsense reasoning on an incomplete

1435 KB. Since no complementary literal for “*Resolves Infection Caused By*” is found, LLM-TH invokes

1436 the LLM to perform a theory resolution step. In this step, it correctly selects “*Disinfects*” as the

1437 clause with the highest priority for theory resolution, enabling it to provide the correct proof.

1438 In this example, all baselines are also able to produce the correct proof. However, the RAG+zero-

1439 shot CoT baseline exhibits minor errors, such as selecting incorrect KB rules in step 4 and neglecting

1440 KB facts in step 5, instead basing its reasoning on an assumption.

F2.2 Geographical QA

1444 **Query:** Can I take a road trip between Ryam Frolovskiy and Biyabanak?

Relevant Rules from the KB:

1446

1447 Facts:

- 1448
- 1449 (1) “*Located in*”(“*Biyabanak*”, “*Iran*”),
- 1450
- 1451 (2) “*Located in*”(“*Ryam Frolovskiy*”, “*Russia*”),
- 1452
- 1453 (3) “*Share Land Borders*”(“*Iran*”, “*Azerbaijan*”)
- 1454
- 1455 (4) “*Share Land Borders*”(“*Azerbaijan*”, “*Russia*”)

1456 Axioms:

1457

- 1458 (5)
 1459
 1460 $\forall A, B, \text{“Share Land Borders”}(A, B) \implies \text{“Land Connected”}(A, B) | A : \text{country}, B : \text{country}$
 1461
 1462 (6)
 1463 $\forall A, B, C, \text{“Land Connected”}(A, B) \wedge \text{“Land Connected”}(B, C)$
 1464 $\implies \text{“Land Connected”}(A, C) | A : \text{location}, B : \text{location}, C : \text{location}$
 1465
 1466 (7)
 1467
 1468 $\forall A, B, C, D, \text{“Located in”}(A, C) \wedge \text{“Located in”}(B, D) \wedge \text{“Land Connected”}(C, D)$
 1469 $\implies \text{“Can Drive Between”}(A, B) | A : \text{location}, B : \text{location}, C : \text{location}, D : \text{location}$
 1470
 1471

Ground Truth Proof:

1472 1- Applying (3) to (5) with the following substitution:
 1473

1474 $\theta = \{A/\text{“Iran”}, B/\text{“Azerbaijan”}\}$ yields:
 1475

1476 (8) $\text{“Land Connected”}(\text{“Iran”}, \text{“Azerbaijan”})$.
 1477

1478 2- Applying (3) to (5) with the following substitution:
 1479

1480 $\theta = \{A/\text{“Azerbaijan”}, B/\text{“Russia”}\}$ yields:
 1481

1482 (9) $\text{“Land Connected”}(\text{“Azerbaijan”}, \text{“Russia”})$.
 1483

1484 3- Applying (8), (9) to (6) with the following substitution:
 1485

1486 $\theta = \{A/\text{“Iran”}, B/\text{“Azerbaijan”}, C/\text{“Russia”}\}$ yields:
 1487

1488 (10) $\text{“Land Connected”}(\text{“Iran”}, \text{“Russia”})$.
 1489

1490 4- Applying (1), (2), (10) to (7) with the following substitution:
 1491

1492 $\theta = \{A/\text{“Biyabanak”}, B/\text{“Ryam Frolovskiy”}, C/\text{“Iran”}, D/\text{“Russia”}\}$ yields:
 1493

1494 (11) $\text{“Can Drive Between”}(\text{“Biyabanak”}, \text{“Ryam Frolovskiy”})$.
 1495

1496 5- Using the commonsense axiom (absent from the KB) that:
 1497

$$\begin{aligned} 1498 & \forall A, B, \text{“Can Drive Between”}(A, B) \\ 1499 & \implies \text{“Can Take Road Trip Between”}(A, B) | A : \text{“location”}, B : \text{“location”} \end{aligned}$$

1500 Therefore, the answer to the query is True.
 1501

LLM-TH:

1502 Proof for query $\text{“Can Drive Between”}(\text{“Biyabanak”}, \text{“Ryam Frolovskiy”})$
 1503

1504 Negated query: $\neg \text{“Can Drive Between”}(\text{“Biyabanak”}, \text{“Ryam Frolovskiy”})$
 1505

1506 Note: In the following sequent calculus formulae, we use “Bi” to refer to “Biyabanak” ,
 1507 “RF” to represent “Ryam Frolovskiy” , “Az” to refer to “Azerbaijan” , “CTRTB” to show
 1508 $\text{“Can Take Road Trip Between”}$, “CDB” to show $\text{“Can Drive Between”}$, “LI” to show “Located In” ,
 1509 “LC” to show “Land Connected” , and “SLB” to show $\text{“Share Land Borders”}$ due to space limi-
 1510 tations.
 1511

Step 1-

$$\frac{\neg \text{“CTRTB”}(\text{“Bi”}, \text{“RF”}) \quad \text{“CDB”}(A, B) \vee \neg \text{“LI”}(A, C) \vee \neg \text{“LI”}(B, D) \vee \neg \text{“LC”}(C, D)}{\neg \text{“LI”}(\text{“Bi”}, C) \vee \neg \text{“LI”}(\text{“RF”}, D) \vee \neg \text{“LC”}(C, D)}$$

Under unification $\theta = \{A/\text{“Biyabanak”}, C/\text{“Ryam Frolovskiy”}\}$.

In this step, by theory resolving \neg “Can Take Road Trip Between”(“Bi”, “RF”) with “Can Drive Between”(A,B), the LLM successfully surfaced the commonsense axiom:

$$\forall A, B, \text{“Can Drive Between”}(A, B) \\ \implies \text{“Can Take Road Trip Between”}(A, B) | A: \text{“location”}, B: \text{“location”}.$$

Step 2-

$$\frac{\neg \text{“LI”}(“Bi”, C) \vee \neg \text{“LI”}(“RF”, D) \vee \neg \text{“LC”}(C, D) \quad \text{“LI”}(“Bi”, “Iran”) \wedge \text{“LI”}(“RF”, “Russia”)}{\neg \text{“LC”}(“Iran”, “Russia”)}$$

Under unification $\theta = \{C/“Iran”, D/“Russia”\}$.

Step 3-

$$\frac{\neg \text{“LC”}(“Iran”, “Russia”) \quad \text{“LC”}(A, C) \vee \neg \text{“LC”}(A, B) \vee \neg \text{“LC”}(B, C)}{\neg \text{“LC”}(“Iran”, B) \vee \neg \text{“LC”}(B, “Russia”)}$$

Under unification $\theta = \{A/“Iran”, C/“Russia”\}$.

Step 4-

$$\frac{\neg \text{“LC”}(“Iran”, B) \quad \text{“LC”}(“Iran”, B) \vee \neg \text{“SLB”}(“Iran”, B)}{\neg \text{“SLB”}(“Iran”, B)}$$

Step 5-

$$\frac{\neg \text{“LC”}(B, “Russia”) \quad \text{“LC”}(B, “Russia”) \vee \neg \text{“SLB”}(B, “Russia”)}{\neg \text{“SLB”}(B, “Russia”)}$$

Step 6-

$$\frac{\neg \text{“SLB”}(“Iran”, B) \vee \neg \text{“SLB”}(B, “Russia”) \quad \text{“SLB”}(“Iran”, “Az”) \wedge \text{“SLB”}(“Az”, “Russia”)}{\perp}$$

Under unification $\theta = \{B/“Az”\}$.

Therefore, the query is proven to be True with the proof score: $\rho_{\text{query}} = (1, 0.994, 6)$

RAG + Zero-shot CoT:

1- From the knowledge base, we have: *PartOf(Biyabanak, Iran)* (from the given KB)

2- *PartOf(Ryam Frolovskiy, Belarus)* (not mentioned in the KB, but assuming it’s true for the sake of the query)

3- Using the rule:

“ $\forall A, B, C, \text{PartOf}(A, B) \wedge \text{PartOf}(B, C) \implies \text{PartOf}(A, C) | A: \text{location}, B: \text{location}, C: \text{location}$ ”,
We can infer: *PartOf(Biyabanak, Asia)* (since Iran is part of Asia)

4- and *PartOf(Ryam Frolovskiy, Europe)* (since Belarus is part of Europe)

5- Using the rule: $\forall A, B, C, \text{LandConnected}(A, B) \wedge \text{LandConnected}(B, C) \implies \text{LandConnected}(A, C) | A: \text{country}, B: \text{country}, C: \text{country}$, we can infer: *Not LandConnected(Iran, Belarus)* (since Iran is in Asia and Belarus is in Europe, and they are not land-connected)

1566 6- Using the rule: “ $\forall A, B, C, D, \text{LocatedIn}(A, C) \wedge \text{LocatedIn}(B, D) \wedge \text{LandConnected}(C, D) \implies$
 1567 $\text{CanDriveBetween}(A, B)$ | $A: \text{location}, B: \text{location}, C: \text{country}, D: \text{country}$ ”, we cannot infer
 1568 $\text{CanDriveBetween}(\text{Biyabanak}, \text{Ryam Frolovskiy})$ since Iran and Belarus are not land-connected.
 1569 Therefore, the answer is **False**.
 1570

1571 RAG + Few-shot CoT:

1572 1- To determine if you can take a road trip between Biyabanak and Ryam Frolovskiy, we need to
 1573 check if Biyabanak and Ryam Frolovskiy are located in countries that are land connected.
 1574 2- From the KB, we know that Biyabanak is part of Iran,
 1575 3- and Ryam Frolovskiy is not mentioned in the KB, but since it is not mentioned,
 1576 we will assume it is in Russia (based on the name).
 1577 4- From the KB, we also know that Iran and Russia are land connected (rule 1).
 1578 5- Using the rule “ $\forall A, B, C, D, \text{LocatedIn}(A, C) \wedge \text{LocatedIn}(B, D) \wedge \text{LandConnected}(C, D) \implies$
 1580 $\text{CanDriveBetween}(A, B)$ | $A: \text{location}, B: \text{location}, C: \text{country}, D: \text{country}$ ”, we can infer that you can
 1581 take a road trip between Biyabanak and Ryam Frolovskiy. Therefore, the answer is **True**.
 1582
 1583

1584 Discussion:

1585 This example further illustrates the successful application of LLM-based theory resolution by LLM-
 1586 TH, where a theory resolution step leverages the LLM’s commonsense knowledge to address KB
 1587 incompleteness. Specifically, in Step 1, the LLM surfaces a commonsense axiom to theory resolve
 1588 $\neg \text{Can Drive Between}$ ” with $\neg \text{Can Take Road Trip Between}$ ”.
 1589 This example also highlights the limitations of the RAG-based baselines. Due to a failure in retriev-
 1590 ing facts about *Ryam Frolovskiy*, the RAG+Zero-shot CoT baseline makes an incorrect assumption
 1591 about its location. Furthermore, subsequent reasoning steps are flawed, resulting in a wrong an-
 1592 swer. Although the RAG+Few-shot CoT baseline provides the correct final answer, its reasoning
 1593 process is not entirely reliable. Notably, because the LLM lacks access to a fact about the location
 1594 of *Ryam Frolovskiy*, it makes an assumption about its location—which happens to be correct in this
 1595 case. However, it also incorrectly references a rule in the KB that states *Iran* and *Russia* are land
 1596 connected, even though no such rule exists in the KB.
 1597 These examples underscore the limited reliability of existing LLM-based baselines when combining
 1598 commonsense reasoning with factual information.
 1599
 1600
 1601
 1602
 1603
 1604
 1605
 1606
 1607
 1608
 1609
 1610
 1611
 1612
 1613
 1614
 1615
 1616
 1617
 1618
 1619